

UNIVERSAL
LIBRARY

OU_164843

UNIVERSAL
LIBRARY

**Compliments of
Scripta Mathematica**

THE SCRIPTA MATHEMATICA LIBRARY
Number Five

Galois Lectures

ADDRESSES DELIVERED BY

JESSE DOUGLAS, PHILIP FRANKLIN

CASSIUS JACKSON KEYSER, LEOPOLD INFELD

AT THE GALOIS INSTITUTE

OF MATHEMATICS

LONG ISLAND UNIVERSITY, BROOKLYN, N. Y.

SCRIPTA MATHEMATICA, *YESHIVA COLLEGE*

NEW YORK, 1941

PRINTED AT THE MORRILL PRESS, FULTON, NEW YORK

CONTENTS

SURVEY OF THE THEORY OF INTEGRATION	
<i>By</i> JESSE DOUGLAS	<i>Page</i> 1
THE FOUR COLOR PROBLEM	
<i>By</i> PHILIP FRANKLIN	<i>Page</i> 49
CHARLES SANDERS PEIRCE AS A PIONEER	
<i>By</i> CASSIUS JACKSON KEYSER	<i>Page</i> 87
THE FOURTH DIMENSION AND RELATIVITY	
<i>By</i> LEOPOLD INFELD	<i>Page</i> 113

SURVEY OF THE THEORY
OF INTEGRATION

By JESSE DOUGLAS

SURVEY OF THE THEORY OF INTEGRATION

CONTENTS

	<i>Page</i>
I. INTRODUCTION	5
1. Area of the circle and the parabola	
2. Area under a continuous curve	
3. Generalization of integral from continuous to arbitrary functions $f(x)$	
II. THE RIEMANN INTEGRAL	11
4. Definition	
5. Criteria for the existence of the Riemann integral	
6. Properties of the Riemann integral	
7. Improper Riemann integrals	
III. THE STIELTJES INTEGRAL	19
8. Definition	
9. Criterion for the existence of the Stieltjes integral, and its properties	
IV. THE LEBESGUE INTEGRAL	25
10. Definition of measure	
11. Properties of measurable sets	
12. Content of a point-set	
13. Measurable functions	
14. Lebesgue integral of a bounded function	
15. Lebesgue integral of an unbounded function	
16. Properties of the Lebesgue integral	
17. Fourier series and transforms	
V. REMARKS ON THE DENJOY INTEGRAL	40
18. Inverse relation between differentiation and integration	

VI. MULTIPLE INTEGRATION	42
19. Riemann multiple integral	
20. Lebesgue multiple integral	
21. Content of a two-dimensional set	
REFERENCES	47

SURVEY OF THE THEORY OF INTEGRATION*

I. INTRODUCTION

1. *Area of the circle and the parabola*

The idea of integration of a function arose from the problem of finding the area bounded by a curve. The ancients could solve this problem only for the case of a circle and a segment of a parabola.

Both these results may be ascribed to Archimèdes. For the circle, he found the value of π in the formula $A = \pi r^2$ to lie between $3\frac{1}{7}$ and $3\frac{10}{71}$. His method, following his predecessors among the Greek geometers, was to express the area of the circle as the limit of the area of an inscribed or a circumscribed regular polygon when the number n of sides was made to increase indefinitely. The close superior and inferior approximations to π given above were obtained by taking $n = 96$.

For the segment of a parabola, Archimedes expressed its area as the sum of an infinite geometric progression of ratio $1/4$, whose n^{th} term is the sum of the areas of 2^n triangles inscribed in the parabolic arc. Fig. 1 shows the triangles at the beginning of the series: $ABC + (ADC + CEB) + \dots$, the law of construction being that the medians such as CL , DM , EN are parallel to the axis of the parabola. The ratio $1/4$ is arrived at by known geo-

*This exposition is intended as an introduction to the theory of integration for non-specialists and students interested in mathematical analysis.

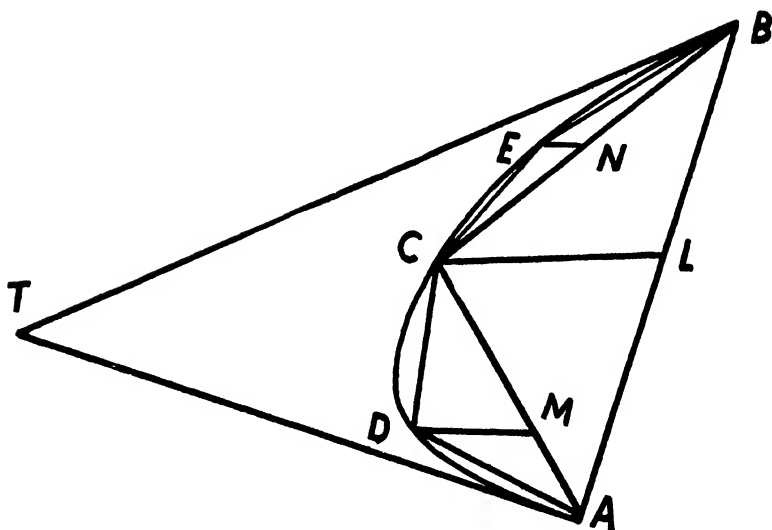


Fig. 1

metric properties of the parabola. By summing the infinite geometric series: $ABC + \frac{1}{4} ABC + \frac{1}{4^2} ABC + \dots$, we obtain with Archimedes the result that the area of the parabolic segment is $\frac{4}{3}$ that of the triangle ABC .

It is seen to be typical of these methods of evaluating area that they depend for their success on the discovery and utilization of some special property of the figure in question. One has to hope in each particular case that he will find some such property by the application of which the required area can be determined. A general method was needed whose systematic application would enable us to evaluate the area of any plane figure whatever, i.e., a figure bounded by one or more closed curves of arbitrary form. The process of integration of a function, as given by the founders of the integral calculus, Newton and Leibniz, and developed and extended by their successors down to

the present day, furnishes exactly the general method desired.

2. Area under a continuous curve

Let $y = f(x)$ be any continuous function defined in the interval ab : $a \leq x \leq b$. Its graph is a curve, represented in fig. 2, which curve together with its two end ordinates Pa , Qb and the segment ab of the x -axis encloses an area,

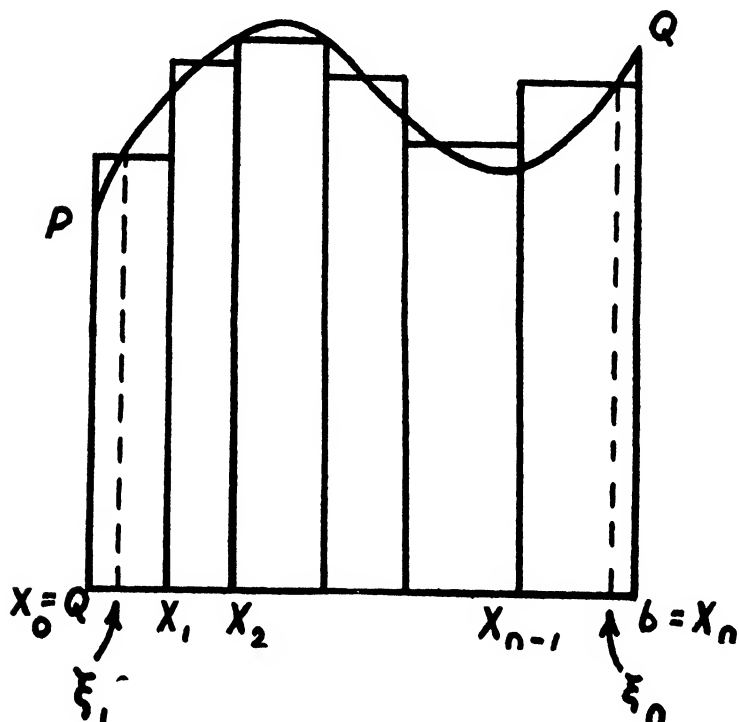


Fig. 2

called the "area under the arc PQ ". Every area bounded by one or more closed curves can be expressed as a sum or difference of areas of this type; for instance, the area

bounded by the curve $PQRSP$ in fig. 3 is equal to: (area under PQ) - (area under RQ) + (area under RS) - (area under PS).

To evaluate the area under the arc PQ , we apply the familiar procedure of the integral calculus; namely, we approximate to the desired area by the sum of a series of rectangles, as in fig. 2. These rectangles are based on any subdivision, or *partition*, of the interval ab by points $x_0 = a$, x_1 , x_2 , . . . , x_{n-1} , $x_n = b$, while the altitude of the rectangle based on the subinterval $x_{k-1}x_k$ is any ordinate of the given curve belonging to a value of x in that subinterval: $f(\xi_k)$, $x_{k-1} \leq \xi_k \leq x_k$. The formula for the total area of the rectangles is

$$(2.1) \quad \sigma = \sum_{k=1}^n f(\xi_k) \delta_k, \quad \text{where} \quad \delta_k = x_k - x_{k-1}.$$

σ is evidently a function of the partition of ab effected by

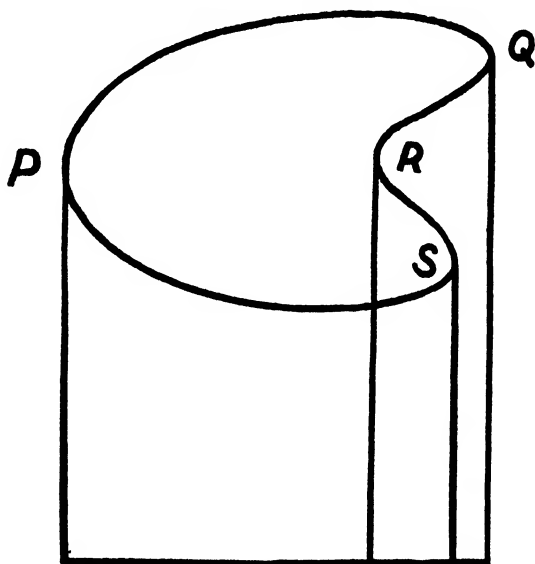


Fig. 3

SURVEY OF THE THEORY OF INTEGRATION

the points x_k and of the choice of the points ξ_k in the sub-intervals $x_{k-1}x_k$.

The function $f(x)$ being continuous, it can be shown by standard methods, involving the *uniform continuity* of $f(x)$ in the interval ab , that when the length of the longest sub-interval δ_k tends to zero,

$$(2.2) \quad \Delta = \max \delta_k \rightarrow 0 \quad (\text{"max" for } k = 1, 2, \dots, n),$$

the value of σ tends to a definite limit I , regardless of the choice of the points ξ_k . This means, precisely, that to every positive number ϵ there corresponds a positive number η such that every partition obeying the condition $\Delta < \eta$ gives a value of σ obeying the condition $|\sigma - I| < \epsilon$.

The limit I thus arrived at is called "the definite integral of $f(x)$ between the limits a and b ," is denoted by the symbol

$$(2.3) \quad I = \int_a^b f(x) dx,$$

and is, *by definition*, the area under the arc PQ of the curve $y = f(x)$. This definition accords with our intuitive ideas concerning area; indeed, it is the only one consistent with the axiom that when one plane region contains another completely, then the area of the whole region is greater than that of its part. To see this, let the points ξ_k first be chosen so that $f(\xi_k)$ is the greatest value of $f(x)$ in each subinterval $x_{k-1}x_k$, and then so that $f(\xi_k)$ is the least such value. Then, according to the stated axiom, the area A under the arc PQ is comprised between the sums σ' , σ'' defined by these two choices of the points ξ_k . Since σ' , σ'' tend to the same limit I , and A is a constant, the conclusion $A = I$ is evidently necessitated.

The connection between the evaluation of I and the determination of an *anti-derivative* or *primitive* of the func-

tion $f(x)$ is well-known, indeed constitutes the fundamental fact of the integral calculus, practice in the application of which forms the substance of college courses in this subject. The relation between definite integral and primitive is: if a continuous function $F(x)$ can be found, defined in the interval ab , such that

$$(2.4) \quad \frac{d}{dx} F(x) = f(x)$$

at every point of this interval, then

$$(2.5) \quad \int_a^b f(x)dx = F(b) - F(a).$$

Conversely, if we regard the upper limit x of the definite integral as variable in the interval ab , the formula

$$(2.6) \quad F(x) = \int_a^x f(x)dx + C,$$

where C is an arbitrary constant, gives the most general primitive of $f(x)$.

3. *Generalization of integral from continuous to arbitrary functions $f(x)$*

We have said that the existence of the limit of the sum σ as $\Delta = \max \delta_k \rightarrow 0$ can be established quite readily if the function $f(x)$ is continuous. It is characteristic of modern function theory to take into consideration not only continuous functions but to allow discontinuities as numerous and of as complicated a nature as possible. The modern viewpoint, in fact, is based on the highly general definition of function given by Dirichlet and Riemann about a hundred years ago, to which they were led largely by their researches in the theory of trigonometric series. Their definition, now established universally, is: " y is a function of x if to every value of x in a certain interval, or other given assemblage of real numbers (or points of the x -axis), there

SURVEY OF THE THEORY OF INTEGRATION

corresponds a definite value of y ." The law of correspondence may be absolutely any whatever, whether defined by an analytic formula of the usual type, with algebraic operations and passages to the limit finite or infinite in number, or by several such formulas, different for different parts of the assemblage (range) of values of x , or in any other way—the only thing essential is the determinate correspondence between x - and y -values.

The question thus naturally arises as to how far and in what way the concepts of integral and primitive, previously given for a continuous function $f(x)$, can be extended to a function in the most general sense of the preceding definition.

The following account is concerned with the various definitions of integral that have been found most useful. These go by the names of their inventors, Riemann, Stieltjes, Lebesgue. Each has its desirable features, useful in applications; each has a high degree of generality, particularly the Lebesgue integral, which applies to a very wide class of functions, including nearly all those ordinarily met with in analysis. Even more general is the definition of integration due to Denjoy, called *totalisation*, which will be referred to briefly.

We shall describe these various types of integral, and discuss their most essential properties. Proofs, however, will be omitted, since these are amply presented in the references listed at the end,¹ which also provide a detailed elaboration of the theory here sketched.

II. THE RIEMANN INTEGRAL

4. *Definition*

The essential idea of Riemann was to form the sum σ of

¹ Citations from this list will be made by means of the author's name.

formula (2.1) for any given function $f(x)$ whatever, and to see if a limit was approached by σ as $\Delta = \max \delta_k \rightarrow 0$. In the affirmative case, the function $f(x)$ is now said to be "integrable in the sense of Riemann," or, briefly, "Riemann integrable," or "integrable R ", and we write

$$(4.1) \quad \int_a^b f(x) dx = \lim_{\Delta \rightarrow 0} \sigma.$$

If the partition $P = (x_0 = a, x_1, x_2, \dots, x_n = b)$ of the interval ab is given, the value of σ still depends on the choice of the intermediate points ξ_k in each subinterval $x_{k-1}x_k$. If M_k, m_k denote respectively the upper and lower bounds of $f(x)$ in $x_{k-1}x_k$, then, by definition of these bounds, $m_k \leq f(\xi_k) \leq M_k$, so that the value of the sum σ is always intermediate between

$$(4.2) \quad S = \sum_{k=1}^n M_k \delta_k, \quad s = \sum_{k=1}^n m_k \delta_k.$$

The possibility of forming the sums S, s rests on the existence of the bounds M_k, m_k as finite quantities; i.e., as a preliminary to defining the Riemann integral, we must assume $f(x)$ to be bounded in ab ;

$$(4.3) \quad m \leq f(x) \leq M \quad (\text{all } x, a \leq x \leq b),$$

where we suppose the finite numbers m, M to be respectively the greatest lower and least upper bounds of $f(x)$ in ab ; i.e., m cannot be increased or M decreased if the relations (4.3) are to continue to subsist.

$f(x)$ being given, S is uniquely determined by the partition P of ab . Since $M_k \geq m, \sum_{k=1}^n \delta_k = b - a$, we have $S \geq m(b - a)$. Thus, the assemblage of values of S formed for all possible partitions P is bounded below; consequently, this assemblage has a *greatest lower bound*, D . It can be proved that D is not only the lower bound of S , but also its *limit*

SURVEY OF THE THEORY OF INTEGRATION

as $\Delta = \max \delta_k \rightarrow 0$; i.e., given $\epsilon > 0$ arbitrarily, a number $\eta > 0$ corresponds such that $|S - D| < \epsilon$ for all partitions P where $\Delta < \eta$. Similar statements apply to s : always $s \leq M(b - a)$; therefore s has a *least upper bound*, d ; and, in fact, d is also the limit of s as $\Delta \rightarrow 0$. It is readily shown that every lower sum s' is less than every upper sum S'' , whether belonging to the same partition of ab or not; accordingly, $d \leq D$.

The numbers D and d were introduced by Darboux in an important memoir on the concept of integral,² and are therefore called respectively the *upper and lower Darboux integrals of $f(x)$* , with the notation:

$$(4.4) \quad D = \overline{\int_a^b} f(x)dx, \quad d = \underline{\int_a^b} f(x)dx.$$

For their existence, only the boundedness of $f(x)$ is required.

As stated, we always have $d \leq D$. When and only when $d = D$, we speak simply of the "integral from a to b of $f(x)$," which we write, as usual, $\int_a^b f(x)dx$. This is the *integral according to the definition of Riemann*. It is the common limit of S and s as $\Delta \rightarrow 0$; and, since $s \leq \sigma \leq S$, it is also the limit of σ as $\Delta \rightarrow 0$, regardless of the choice of the points ξ_k :

$$(4.5) \quad \int_a^b f(x)dx = \lim S = \lim s = \lim \sigma \text{ as } \Delta \rightarrow 0.$$

5. Criteria for the existence of the Riemann integral

This definition raises the question: what, exactly, is the class of functions $f(x)$ for which the limit required in the

² *Sur les fonctions discontinues*, Annales de l'École Normale Supérieure, 1875.

definition of Riemann integration exists? As has been remarked, every continuous function is Riemann integrable. It is easy to show that every function bounded in an interval ab , and with only a finite number of points of discontinuity, is integrable R . It is, further, true that a denumerable infinity of discontinuities may be allowed to a bounded function, and this function will still be integrable R .

Certainly, however, not every bounded function is integrable R . Consider, for instance, the function $\varphi(x)$, due to Dirichlet, equal to 1 when x is rational, and equal to 0 when x is irrational. Obviously, $\varphi(x)$ is everywhere discontinuous. Then we evidently have

$$(5.1) \quad \int_0^1 \varphi(x) dx = 1, \quad \int_0^1 \varphi(x) dx = 0;$$

accordingly, $\varphi(x)$ is not integrable in the sense of Riemann.

This example shows that a function may be so radically discontinuous as not to be capable of Riemann integration. What kind and degree of discontinuity is compatible with integrability R ? The neatest form of answer was given by Lebesgue: *a necessary and sufficient condition that a bounded function $f(x)$, $a \leq x \leq b$, be Riemann integrable is that the points of discontinuity of $f(x)$ form a set E of zero measure.*

"Zero measure" means that the points of E can be covered by a denumerable infinity of intervals I_n whose total length is arbitrarily small: $\sum_{n=1}^{\infty} I_n < \epsilon$, where ϵ denotes any

preassigned positive number. Any set E finite in the number of its points, or denumerably infinite, is easily seen to be of zero measure.^{2a} But a set E may contain a non-denumerable infinity of points, and still be of zero measure—such, for example, is the well-known set due to

^{2a} See §10, last paragraph.

SURVEY OF THE THEORY OF INTEGRATION

Cantor. This consists of all real numbers in the interval 01 which can be expressed in the radix-system based on 3 as follows:

$$(5.2) \quad x = \frac{a_1}{3} + \frac{a_2}{3^2} + \dots + \frac{a_n}{3^n} + \dots,$$

where the digits $a_1, a_2, \dots, a_n, \dots$ are restricted to the values 0, 2, the digit 1 being excluded. The measure of a set expresses its metrical extent, not the cardinal number of its points.

Another form of necessary and sufficient condition for the Riemann integrability of a function $f(x)$ is the following: *the mean oscillation Ω of $f(x)$ shall be zero*. The oscillation of $f(x)$ in an interval δ_k is, by definition, $\omega_k = M_k - m_k$. The mean oscillation of $f(x)$ relative to any given partition of its interval of definition ab is

$$(5.3) \quad \omega = \frac{\sum_{k=1}^n \omega_k \delta_k}{b-a}$$

Ω is defined to be the greatest lower bound of ω for all possible partitions of ab . The proof of the stated condition follows almost immediately from the definition of Riemann integral.

6. Properties of the Riemann integral

By either of the preceding criteria, it is readily seen that if $f(x)$, $g(x)$ are any two Riemann integrable functions, then $f(x) + g(x)$, $f(x) \cdot g(x)$, $f(x)/g(x)$ are likewise Riemann integrable, provided, in the case of division, that $g(x)$ is nowhere equal to 0 and $1/g(x)$ is bounded. Furthermore, the Riemann integral is additive, both with respect to the functions integrated and as to the interval of integration; i.e.:

$$(6.1) \quad \int_a^b [f(x) + g(x)]dx = \int_a^b f(x)dx + \int_a^b g(x)dx,$$

$$(6.2) \quad \int_a^b f(x)dx + \int_b^c f(x)dx = \int_a^c f(x)dx.$$

Again, the Riemann integral obeys the "law of the mean":

$$(6.3) \quad \int_a^b f(x)dx = \mu(b - a), \text{ where } m \leq \mu \leq M,$$

m, M being the bounds of $f(x)$ in ab .

As to the relation of Riemann integration to the limit process, we have the theorem: *if a sequence of Riemann integrable functions, $f_n(x)$, $a \leq x \leq b$, tends uniformly to a limit function $f(x)$, then $f(x)$ is Riemann integrable and*

$$(6.4) \quad \int_a^b f(x)dx = \lim_{n \rightarrow \infty} \int_a^b f_n(x)dx.$$

The stipulation "uniformly" is essential, since the mere limit of functions integrable R need not be integrable R . For instance, let the rational numbers in the interval 01 be arranged in a sequence r_n ($n = 1, 2, 3, \dots$), as is known to be possible. Let $f_n(x)$ be equal to zero in the interval 01 except at the points r_1, r_2, \dots, r_n , where its value shall be 1. Then $\lim_{n \rightarrow \infty} f_n(x) = \varphi(x)$, where $\varphi(x)$ is the

function previously mentioned, equal to 1 when x is rational, equal to 0 when x is irrational. As has been observed, $\varphi(x)$ is not integrable R , although, clearly, each $f_n(x)$ is so integrable, indeed,

$$\int_0^1 f_n(x)dx = 0 \text{ for } n = 1, 2, 3, \dots$$

In practical applications (engineering), Riemann integration is the kind that is actually used; for instance, in obtaining an area or a volume by the trapezoidal rule or by Simpson's rule. The approximate value of the integral

furnished by these rules is essentially of the type of the approximating sums that characterize the Riemann integral.

7. Improper Riemann integrals

Let the function $f(x)$ be defined in the interval ab , bounded and Riemann integrable in every subinterval $(a + \epsilon, b)$ where $\epsilon > 0$, but unbounded in the entire interval ab , i.e., becoming unbounded when $x \rightarrow a$. Then, following Cauchy, we may define the "improper integral":

$$(7.1) \quad \lim_{\epsilon \rightarrow 0} \int_{a+\epsilon}^b f(x) dx = \int_a^b f(x) dx,$$

provided that the indicated limit exists. A similar definition applies if the singular nature of $f(x)$ is at the other end b of the interval of integration. If $f(x)$ is singular at both end-points a, b , we choose any intermediate point c and define:

$$(7.1') \quad \int_a^b f(x) dx = \int_a^c f(x) dx + \int_c^b f(x) dx,$$

in case each improper integral on the right exists; the fact that this definition does not depend on the position of the point c is obvious.

For example, the function

$$(7.2) \quad f(x) = \frac{A}{x^n} \quad (n > 0)$$

has the properties stated above in the interval 01 (let $f(0)$ have any chosen value, by definition). We have

$$(7.3) \quad \int_{\epsilon}^1 \frac{A}{x^n} dx = \frac{A}{1-n} (1 - \epsilon^{1-n}), \text{ if } n \neq 1,$$

$$(7.4) \quad \int_{\epsilon}^1 \frac{A}{x} dx = -A \log \epsilon.$$

Accordingly, if $0 < n < 1$, we have, by the definition (7.1),

the existence and value of the following improper integral:

$$(7.5) \quad \int_0^1 \frac{A}{x^n} dx = \frac{A}{1-n}.$$

On the other hand, if $n \geq 1$, then this improper integral does not exist as a finite value.

More generally, if $f(x)$ can be expressed in the form

$$(7.6) \quad f(x) = \frac{g(x)}{x^n},$$

where $0 < n < 1$, and $g(x)$ stays bounded in the interval 01 and is Riemann integrable in every subinterval $\epsilon 1$, then $\int_0^1 f(x)dx$ exists as an improper Riemann integral. On the contrary, if

$$(7.7) \quad f(x) > \frac{A}{x^n},$$

where A is a positive constant and $n \geq 1$, then the improper integral $\int_0^1 f(x)dx$ does not exist.

If there are only a finite number of points of unboundedness of $f(x)$ in the interval ab , namely, c_1, c_2, \dots, c_k , then we define

$$(7.8) \quad \int_a^b f(x)dx = \int_a^{c_1} + \int_{c_1}^{c_2} + \dots + \int_{c_{k-1}}^{c_k} + \int_{c_k}^b,$$

where each integral on the right is improper, to be defined as a limit by the formula (7.1) or (7.1').

The notion of improper integral may also be extended to the case where the derived set E' , meaning the set of limit-points of the set E of singularities of $f(x)$, consists of only a finite number of points. More generally, let $E^{(\alpha)}$ denote the derived set of E of any finite or transfinite order α .³ If $E^{(\alpha)}$ consists of only a finite number of points,

³ See Lebesgue, pp.12-14, 320-324.

then $\int_a^b f(x)dx$ can be defined as an improper integral by successive limiting processes.

If the interval of integration is infinite at one end, $a \infty$, we define

$$(7.9) \quad \int_a^\infty f(x)dx = \lim_{b \rightarrow \infty} \int_a^b f(x)dx,$$

provided that the indicated limit exists. Similarly, for a function $f(x)$ defined over the entire real axis, we have the definition

$$(7.10) \quad \int_{-\infty}^{+\infty} f(x)dx = \lim_{\substack{a \rightarrow -\infty \\ b \rightarrow +\infty}} \int_a^b f(x)dx.$$

Thus, if $n > 0$ and $\neq 1$,

$$(7.11) \quad \int_1^\infty \frac{A}{x^n} dx = \lim_{b \rightarrow \infty} \frac{A}{n-1} \left(1 - \frac{1}{b^{n-1}} \right) = \frac{A}{n-1} \text{ if } n > 1,$$

and does not exist if $n < 1$. When $n = 1$, we have

$$(7.12) \quad \int_1^\infty \frac{A}{x} dx = \lim_{b \rightarrow \infty} A \log b,$$

which does not exist as a finite value.

III. THE STIELTJES INTEGRAL

8. Definition

Stieltjes integration, introduced by Thomas Jan Stieltjes in a memoir of 1894,⁴ is the integration of $f(x)$ with respect to a "determining function" $\alpha(x)$, instead of merely with respect to the independent variable x .

Let $\alpha(x)$ denote, to start with, a *monotonic increasing* continuous function, and let $f(x)$ also be continuous in its

⁴ On the subject of continued fractions, published in the *Annales de la Faculté des Sciences de Toulouse*.

interval ab of definition. Relative to any partition ($x_0 = a$, $x_1, x_2, \dots, x_n = b$) of ab , let us form the two sums (upper and lower)

$$(8.1) \quad S = \sum_{k=1}^n M_k[\alpha(x_k) - \alpha(x_{k-1})],$$

$$s = \sum_{k=1}^n m_k[\alpha(x_k) - \alpha(x_{k-1})].$$

We have, for their difference

$$(8.2) \quad S - s = \sum_{k=1}^n \omega_k[\alpha(x_k) - \alpha(x_{k-1})] \leq \omega[\alpha(b) - \alpha(a)],$$

where ω denotes the greatest of the oscillations ω_k of $f(x)$ in the intervals $x_{k-1}x_k$ of the given partition. In the derivation of the inequality (8.2), the assumed monotonic increasing nature of $\alpha(x)$ intervenes in the use of the condition that each bracket is non-negative.

By a standard theorem on continuous functions $f(x)$, (uniform continuity),

$$(8.3) \quad \omega = \max_k \omega_k \rightarrow 0 \text{ as } \Delta = \max_k (x_k - x_{k-1}) \rightarrow 0;$$

therefore $S - s \rightarrow 0$ as $\Delta \rightarrow 0$. Since always $S \geq s$ —and indeed, as can be shown, whether the upper and lower sums S, s belong to the same or to different partitions of ab —it follows quite readily that as $\Delta \rightarrow 0$, S and s tend to the same limit. This limit is called “the Stieltjes integral of $f(x)$ with respect to $\alpha(x)$,” and is denoted by

$$(8.4) \quad \int_a^b f(x) d\alpha(x) = \lim_{\Delta \rightarrow 0} \sum_{k=1}^n f(\xi_k)[\alpha(x_k) - \alpha(x_{k-1})],$$

$$(x_{k-1} \leq \xi_k \leq x_k),$$

the equation being justified, regardless of the choice of intermediate points ξ_k , because the value of the sum on the right is comprised between S and s .

SURVEY OF THE THEORY OF INTEGRATION

The continuous monotonic increasing function $y = \alpha(x)$ always has an inverse of the same nature,⁵ $x = \beta(y)$, and it is evident, by comparing (8.4) with (4.5), (2.1), that the Stieltjes integral can also be expressed as a Riemann integral:

$$(8.5) \quad \int_a^b f(x) d\alpha(x) = \int_{\alpha(a)}^{\alpha(b)} f(\beta(y)) dy.$$

Obviously, however, the continuity of the monotonic function $\alpha(x)$ intervened nowhere in the reasoning associated with the definition (8.4); we may therefore allow $\alpha(x)$ in that definition to have any discontinuities possible for a bounded monotonic function. It is well-known that such points of discontinuity c are at most denumerably infinite in number, and that each is of the nature of unequal one-side limits: $\alpha(c+0) = \lim_{h \rightarrow 0} \alpha(c+h)$, $\alpha(c-0) = \lim_{h \rightarrow 0} \alpha(c-h)$, $h > 0$, with a "saltus" $\sigma(c) = \alpha(c+0) - \alpha(c-0)$.

If the interval ab is divided into m subintervals by division-points c_1, c_2, \dots, c_{m-1} , and in each of the intervals $ac_1, c_1c_2, \dots, c_{m-2}c_{m-1}, c_{m-1}b$, the function $\alpha(x)$ has a constant value: $\alpha_1, \alpha_2, \dots, \alpha_{m-1}, \alpha_m$, respectively, then

$$(8.6) \quad \int_a^b f(x) d\alpha(x) = \sum_{k=1}^{m-1} \sigma_k f(c_k),$$

where σ_k denotes the saltus of $\alpha(x)$ at c_k : $\sigma_k = \alpha_{k+1} - \alpha_k$. This shows how a finite sum may be included under the notion of Stieltjes integral; evidently, it is nowise necessary for the validity of (8.6) that the piece-wise constant function $\alpha(x)$ be monotonic, so that the coefficients σ_k may be any positive or negative numbers.

⁵ In this statement, we suppose $\alpha(x)$ to be *strictly* increasing, i.e., if $x_1 < x_2$, then $\alpha(x_1) < \alpha(x_2)$ —not merely $\alpha(x_1) \leq \alpha(x_2)$.

The Stieltjes integral may easily be extended to functions $\alpha(x)$ expressible as the difference of two monotonic increasing functions:

$$(8.7) \quad \alpha(x) = \beta(x) - \gamma(x),$$

by observing that in this case also the limit

$$(8.8) \quad \int_a^b f(x) d\alpha(x) = \lim_{\Delta \rightarrow 0} \left\{ \sum_{k=1}^n f(\xi_k) [\beta(x_k) - \beta(x_{k-1})] \right. \\ \left. - \sum_{k=1}^n f(\xi_k) [\gamma(x_k) - \gamma(x_{k-1})] \right\}$$

must exist, indeed, is equal to $\int_a^b f(x) d\beta(x) - \int_a^b f(x) d\gamma(x)$.

How can we characterize functions $\alpha(x)$ expressible in the form (8.7) with $\beta(x)$, $\gamma(x)$ monotonic increasing? Let $P = (x_0 = a, x_1, x_2, \dots, x_{n-1}, x_n = b)$ denote any partition of ab ; then

$$(8.9) \quad \alpha(x_k) - \alpha(x_{k-1}) = [\beta(x_k) - \beta(x_{k-1})] \\ - [\gamma(x_k) - \gamma(x_{k-1})];$$

consequently,

$$(8.10) \quad |\alpha(x_k) - \alpha(x_{k-1})| \leq |\beta(x_k) - \beta(x_{k-1})| \\ + |\gamma(x_k) - \gamma(x_{k-1})|,$$

$$(8.11) \quad \sum_{k=1}^n |\alpha(x_k) - \alpha(x_{k-1})| \leq [\beta(b) - \beta(a)] \\ + [\gamma(b) - \gamma(a)].$$

The sum on the left, called the "variation of $\alpha(x)$ relative to the partition P ," is therefore bounded for all possible partitions of ab ; consequently, it has a finite least upper bound V , called the "variation of $\alpha(x)$ in the interval ab ". Accordingly, a necessary condition for the representability of $\alpha(x)$ as the difference of two monotonic increasing functions is that $\alpha(x)$ be of *bounded variation*, meaning the finiteness of V .

It is proved in text-books on the theory of functions of a real variable⁶ that this condition is also sufficient. It follows that the limit (8.8) called for in the definition (8.4) of the Stieltjes integral certainly exists whenever the integrand function $f(x)$ is continuous and the determining function $\alpha(x)$ is of bounded variation.

9. *Criterion for the existence of the Stieltjes integral, and its properties*

If $\alpha(x)$ denote a given function of bounded variation, a necessary and sufficient condition on $f(x)$ for the existence of the Stieltjes integral $\int_a^b f(x) d\alpha(x)$ is that the set E of points of discontinuity of $f(x)$ have an “ α -measure” equal to zero. The definition of “ α -measure,” a generalization of ordinary measure, to which it reduces when $\alpha(x) = x$, is given in the book of Lebesgue⁷ cited at the end. If $\alpha(x)$ is a monotonic increasing function, the condition of zero α -measure signifies the following property of the set E : it is possible to cover E with a system of intervals $J = \{a_n b_n\}$ finite or denumerably infinite in number, such that

$$(9.1) \quad \sum_n [\alpha(b_n) - \alpha(a_n)] < \epsilon,$$

where ϵ denotes an arbitrarily preassigned positive number. The “covering” of E by J means that, for each point x of E , there exists an interval of the given system J which contains the point x *strictly in its interior*. This strictly interior position of the point x is necessitated by possible discontinuities of $\alpha(x)$ —by contrast, when $\alpha(x) = x$, the point x may also be an extremity of its containing interval.

⁶ E.g., Lebesgue, p.52.

⁷ Pages 276–281.

One should be careful to observe that the condition of only a finite number k of points of discontinuity of $f(x)$ does not guarantee the existence of the Stieltjes integral

$\int_a^b f(x) d\alpha(x)$. Supposing, for simplicity that $\alpha(x)$ is monotonic increasing, it is further necessary (and sufficient) that all of these points of discontinuity of $f(x)$ be points of continuity of $\alpha(x)$; otherwise, the finite set of k points has a positive α -measure, equal to the sum of the values of the saltus of $\alpha(x)$, namely $\alpha(c+0) - \alpha(c-0)$, extended over the simultaneous points of discontinuity of $f(x)$ and $\alpha(x)$.

An important formula in the theory of Stieltjes integration is that of *integration by parts*, namely:

$$(9.2) \quad \int_a^b f(x) d\alpha(x) = [f(b)\alpha(b) - f(a)\alpha(a)] - \int_a^b \alpha(x) df(x).$$

The interpretation is that the existence of either of the two integrals which appear implies that of the other, and that the two are connected by the stated formula. A consequence is that $\int_a^b f(x) d\alpha(x)$ will exist if the integrand $f(x)$ is of bounded variation and the determining function $\alpha(x)$ is continuous.

The Stieltjes integral is important in physical applications. Suppose, as a typical illustration, that the interval ab of the x -axis is loaded with matter so that the mass in the interval $0x$ is expressed by the monotonic increasing function $\alpha(x)$. A positive mass m concentrated at the point c will be represented by a discontinuity of $\alpha(x)$ at that point with a saltus $\alpha(c+0) - \alpha(c-0) = m$. Then the *moment* of the total mass about the origin is expressed by

the Stieltjes integral $\int_a^b x d\alpha(x)$, and the moment of n^{th} order by $\int_a^b x^n d\alpha(x)$.

IV. THE LEBESGUE INTEGRAL

10. *Definition of measure*

The important generalization of the idea of integral introduced by Lebesgue about 1902 is based on the notion of *measure* of a set of points E .

Throughout, we shall regard all our point-sets E as belonging to a given interval ab . As defined in the previous section, a *covering-system of intervals*, $J = \{I_n = a_n b_n\}$, of the set E is a system of intervals, finite or denumerably infinite in number, overlapping or not, such that each point x of E is contained (whether as interior or end-point) in some interval $I_{n(x)}$ of the system. We denote by $l(J)$ the total length, finite or infinite, of the intervals which compose the system J :

$$(10.1) \quad l(J) = \sum_n (b_n - a_n).$$

The lower bound of $l(J)$ for all possible covering systems of the set E is termed the *outer measure of E* , notationally:

$$(10.2) \quad \overline{m}(E) = \underline{B}l(J) \quad [\text{all } J \text{ covering } E].$$

Let CE denote the "complementary set" or *complement* of E with respect to the interval ab ; i.e., CE consists exactly of those points of ab which do not belong to E . Then we define the *inner measure* of E by the formula

$$(10.2') \quad \underline{m}(E) = (b - a) - \overline{m}(CE).$$

It is readily shown that, for every set E , we have

$$(10.3) \quad \underline{m}(E) \leq \overline{m}(E).$$

The vital element in the proof is the well-known Heine-Borel theorem: *if all the points of a closed interval ab are covered by a given system of intervals, $J = \{I_n\}$, $n = 1, 2, 3, \dots$, then there is a finite subset of these intervals, $j = \{I_m\}$, $m = n_1, n_2, \dots, n_k$, such that j suffices to cover the interval ab .* This theorem also plays an important role at many other stages in the theory of point-sets and of measure.

In case the $=$ sign really holds in (10.3), i.e., inner measure is equal to outer measure, we say that the set E is *measurable* and define its *measure* $m(E)$ to be the common value of inner and outer measure:

$$(10.4) \quad m(E) = \underline{m}(E) = \overline{m}(E).$$

Evidently, the complement of any measurable set is measurable, and we have

$$(10.5) \quad m(CE) = (b - a) - m(E).$$

It is easily seen from these definitions that any interval is measurable, with a measure equal to its length, and that any sum of a finite number of non-overlapping intervals is measurable, with a measure equal to the total length of the component intervals. Any finite number of points is immediately seen to have the measure zero. But, further, *any denumerable infinity of points, $(x_1, x_2, \dots, x_n, \dots)$, is of measure zero.* For, $\epsilon > 0$ being assigned arbitrarily, we may express it as the sum of an infinite series:

$$(10.6) \quad \epsilon = \frac{\epsilon}{2} + \frac{\epsilon}{2^2} + \dots + \frac{\epsilon}{2^n} + \dots,$$

and cover the point x_n with an interval of length $\frac{\epsilon}{2^n}$, for instance, the interval $\left(x_n - \frac{1}{2} \cdot \frac{\epsilon}{2^n}, x_n + \frac{1}{2} \cdot \frac{\epsilon}{2^n}\right)$. This ex-

presses the definition of zero measure for the given set $(x_1, x_2, \dots, x_n, \dots)$.

11. Properties of measurable sets

The *logical sum*, or simply the *sum*, of any finite number or denumerable infinity of sets E_n , $n = 1, 2, 3, \dots$, means the set E which consists exactly of those points x *belonging to at least one of the sets* E_n . We write, in the case of a finite number m of sets:

$$(11.1) \quad E = E_1 + E_2 + \dots + E_m,$$

and, in the case of a denumerable infinity:

$$(11.2) \quad E = E_1 + E_2 + \dots + E_n + \dots$$

The *logical product*, or simply the *product*, of any finite or denumerably infinite number of sets E_n means the set E which consists exactly of those points x *belonging to all of the sets* E_n ; we write

$$(11.3) \quad E = E_1 \cdot E_2 \cdot \dots \cdot E_m,$$

or

$$(11.4) \quad E = E_1 \cdot E_2 \cdot \dots \cdot E_n \cdot \dots,$$

according as there are a finite number or a denumerable infinity of sets E_n .

An important theorem on measure is that *the sum and the product of any finite number or denumerable infinity of measurable sets is itself measurable*. Further, if the sets E_n are mutually exclusive (have no points in common), then we have the formula

$$(11.5) \quad m(E_1 + E_2 + \dots + E_m) = m(E_1) + m(E_2) + \dots + m(E_m),$$

and also

$$(11.6) \quad m(E_1 + E_2 + \dots + E_n + \dots) = m(E_1) \\ + m(E_2) + \dots + m(E_n) + \dots,$$

the last being a convergent infinite series. In words: *the measure of the sum of any finite or denumerably infinite number of mutually exclusive sets is equal to the sum of their measures.*

If the set F is part of the set E , then those points of E which do not belong to F form a set, called the (logical) difference of E and F , and denoted by $E - F$:

$$(11.7) \quad E - F = E \cdot CF.$$

Let E and F be measurable; then CF is measurable; therefore $E - F$ is measurable. Since $(E - F) + F = E$, and $E - F$ and F are mutually exclusive, we have, by applying (11.5),

$$(11.8) \quad m(E - F) = m(E) - m(F).$$

Thus, the difference of two measurable sets is measurable.

By repeated application of the preceding theorems, we derive a very large category of measurable sets, and, at the same time, their measures. The measure of an interval (its length) and of a point (zero) is known. Starting with points and intervals, we may apply the operations of forming (1) the logical sum, (2) the logical product, (3) the logical difference, (4) the complement, any finite or denumerably infinite number of times in succession; each application is to sets already shown to be measurable, and whose measures are known. All sets which can be arrived at in this way may be called *constructible*; such are all the sets practically used in analysis. Because of the particular consideration of these sets by the mathematician E. Borel, such sets are said to be "Borel measurable," or "measurable B ." The more general definition of measure associated with

the formulas (10.2), (10.2'), (10.4) is due to Lebesgue. All sets measurable B are measurable according to Lebesgue, and with the same measure, but sets exist which are measurable according to Lebesgue yet not measurable B . The Lebesgue definition is thus wider in its scope, though the Borel definition includes most of the important sets actually considered in analysis.

Are all point-sets measurable in the sense of Lebesgue? The only point-sets not Lebesgue measurable which have been constructed involve in their formation the so-called "axiom of choice," or "Zermelo's axiom," one of whose formulations is the following: "If $\{C\}$ is an infinite class of classes C , mutually exclusive, then a class D exists which consists of exactly one element from each class C ". The logical admissibility of this axiom has been the subject of much controversy among mathematicians and logicians in recent years.⁸

12. Content of a point-set

It is helpful in understanding both ideas to contrast with the notion of measure of a point-set that of its *content*.

If, in the preceding definition of measure, we restrict the covering-system of the set E to contain only a *finite* number of intervals:

$$(12.1) \quad j = \{I_n = a_n b_n\} \quad (n = 1, 2, \dots, m),$$

then similar formulas give us the definition of the "outer content" $\bar{c}(E)$, the "inner content" $\underline{c}(E)$, and, in the case of equality: $\bar{c}(E) = \underline{c}(E)$, the "content" $c(E)$. These formulas are, precisely,

$$(12.2) \quad \bar{c}(E) = \underline{B}l(j) = \underline{B} \sum_{n=1}^m (b_n - a_n)$$

⁸ See Lebesgue, p.114, footnote 1.

for all finite covering-systems j of E ,

$$(12.3) \quad \underline{c}(E) = (b - a) - \bar{c}(CE),$$

$$(12.4) \quad c(E) = \bar{c}(E) = \underline{c}(E) \text{ in case } \bar{c}(E) = \underline{c}(E)$$

—always, it should be remarked, $\bar{c}(E) \geq \underline{c}(E)$.

The notion of content, older than that of measure, is thus seen to be distinguished from the more modern and useful notion by the restriction of the covering-system j to a finite number of intervals, whereas the covering-system J in the definition of measure is permitted to contain a denumerable infinity of intervals. A set that has a content, $c(E) = \bar{c}(E) = \underline{c}(E)$, is sometimes said to be “measurable J ,” where this letter recalls the name of C. Jordan, who first introduced the notion. We may notice that $\bar{c}(E) \geq \bar{m}(E)$, since the class of covering-systems j is part of the class of covering-systems J . It follows that $\bar{c}(E) \geq \bar{m}(E) \geq \underline{m}(E) \geq \underline{c}(E)$; consequently, every set which has a content also has a measure, which is equal to its content.

The converse is not true, as may be seen by the following example. Let E denote the set of rational points of the interval 01. Then, evidently, we have

$$(12.5) \quad \bar{c}(E) = 1, \quad \underline{c}(E) = 0,$$

so that the set has no definite content. On the other hand, the measure of the same set is zero, on account of the denumerability of the rational numbers.

By the *characteristic function* of a set E , we mean the function $\psi(x)$ equal to 1 when x belongs to E , equal to 0 when x does not belong to E . We easily see that the content of any given set E is the Riemann integral of its characteristic function; in fact:

$$(12.6) \quad \bar{c}(E) = \overline{\int \psi(x) dx}, \quad \underline{c}(E) = \underline{\int \psi(x) dx},$$

$$c(E) = \int_{(R)} \psi(x) dx.$$

The measure of the set E , on the other hand, is equal to the Lebesgue integral of its characteristic function:

$$(12.7) \quad m(E) = \int_{(L)} \psi(x) dx,$$

as will be obvious when the definition of Lebesgue integral is given in §14.

13. Measurable functions

A function $f(x)$, defined on an interval ab or on a measurable set E , is said to be *measurable* if and only if, for every given real number A , the set of points

$$(13.1) \quad E\{f(x) > A\},$$

i.e., the set of points x where the condition $f(x) > A$ is obeyed, is measurable. This implies that, in a similar notation, each of the following sets is measurable for all values of A :

$$(13.2) \quad E\{f(x) \geq A\}, \quad E\{f(x) < A\}, \quad E\{f(x) \leq A\}.$$

Conversely, the measurability for every A of any of these sets implies the measurability of the set (13.1) for every A . Also, if $f(x)$ is measurable, the set $E\{f(x) = A\}$ is measurable for every value of A .

The sum of any finite number of measurable functions is measurable; the same is true for the product, and for the quotient in so far as the denominator is nowhere equal to zero. If $f_n(x)$, $n = 1, 2, 3, \dots$, denote any infinite sequence of measurable functions which tend to a limit function $f(x)$, then $f(x)$ is measurable—the convergence need not be uniform or of any other special kind. More

generally, the inferior and superior limits of any sequence of measurable functions $f_n(x)$:

$$\liminf f_n(x), \quad \limsup f_n(x),$$

are measurable. Here, by definition,

$$(13.3) \quad \liminf f_n(x) = \lim_{n \rightarrow \infty} \lim_{m \rightarrow \infty} \text{least } \{f_n(x), \\ f_{n+1}(x), \dots, f_{n+m}(x)\},$$

$$(13.4) \quad \limsup f_n(x) = \lim_{n \rightarrow \infty} \lim_{m \rightarrow \infty} \text{greatest } \{f_n(x), \\ f_{n+1}(x), \dots, f_{n+m}(x)\}.$$

Briefly summarized, all the usual processes of analysis applied to measurable functions give measurable functions as a result.

14. *Lebesgue integral of a bounded function*

Let $f(x)$ denote any measurable function defined on a measurable set E , and let $f(x)$ be bounded:

$$(14.1) \quad b \leq f(x) \leq B,$$

where b, B denote the closest bounds of $f(x)$.

Two features are typical of the Lebesgue definition of integral, as contrasted with that of Riemann: (1) partitions are made in the interval bB of the functional values $f(x)$, instead of in the interval ab of the argument values x ;⁹ (2) the notion of measure of a point-set is applied, instead of the length of an interval or the content of a point-set.

Specifically, let $P = (b_0 = b, b_1, b_2, \dots, b_{n-1}, b_n = B)$ denote any partition of the interval bB . We form the

⁹ The slight imperfection of notation due to the double occurrence of the letter b in the intervals ab, bB should not cause the least ambiguity.

“upper and lower sums,” S and s , of $f(x)$ relative to P as follows:

$$(14.2) \quad S = b \cdot mE\{f(x) = b\} + \sum_{k=1}^n b_k \cdot mE\{b_{k-1} < f(x) \leq b_k\},$$

$$(14.3) \quad s = \sum_{k=1}^n b_{k-1} \cdot mE\{b_{k-1} \leq f(x) < b_k\} + B \cdot mE\{f(x) = B\}.$$

Let λ denote the length of the greatest of the subintervals $b_k - b_{k-1}$ of the partition P . Then it can be proved that, as $\lambda \rightarrow 0$, the sums S , s always tend to the same limit L . This limit is called the Lebesgue integral of $f(x)$ over the set E :

$$(14.4) \quad L = \lim_{\lambda \rightarrow 0} S = \lim_{\lambda \rightarrow 0} s = \int_E f(x) dx.$$

Thus, every measurable function that is bounded is integrable in the sense of Lebesgue. In contrast, not every bounded function is integrable according to Riemann.

15. *Lebesgue integral of an unbounded function*

Let $f(x)$ denote an unbounded measurable function defined on any given measurable set E . We define $f_1(x)$, $f_2(x)$ —respectively, the positive and negative parts of $f(x)$ —as follows:

$$(15.1) \quad \begin{aligned} f_1(x) &= f(x) \text{ if } f(x) \geq 0, f_1(x) = 0 \text{ if } f(x) < 0; \\ f_2(x) &= 0 \text{ if } f(x) \geq 0, f_2(x) = -f(x) \text{ if } f(x) < 0. \end{aligned}$$

Of course, $f_1(x) \geq 0$, $f_2(x) \geq 0$ for all x in E , and

$$(15.2) \quad f(x) = f_1(x) - f_2(x).$$

Thus it suffices to consider unbounded positive-valued functions $f_1(x)$, $f_2(x)$; we shall then define:

$$(15.3) \quad \int_E f(x)dx = \int_E f_1(x)dx - \int_E f_2(x)dx.$$

Supposing, then, $f(x) \geq 0$, we define, for any $M > 0$,

$$(15.4) \quad f_M(x) = \begin{cases} f(x) & \text{if } f(x) \leq M, \\ M & \text{if } f(x) > M, \end{cases} = \text{smaller } \{f(x), M\}.$$

$f_M(x)$ may be called " $f(x)$ truncated to the value M ".

$f_M(x)$ being bounded and measurable, the Lebesgue integral $\int_E f_M(x)dx$ is defined by the preceding section.

Evidently, it is an increasing function of M , and therefore has a finite or infinite limit as $M \rightarrow \infty$. In case this limit is finite, the function $f(x)$ is called *summable*, and we define

$$(15.5) \quad \int_E f(x)dx = \lim_{M \rightarrow \infty} \int_E f_M(x)dx.$$

An unbounded function $f(x)$ whose values are capable of either sign is called summable if the related functions $f_1(x)$, $f_2(x)$ are both summable, and then the Lebesgue integral of $f(x)$ is defined, as has been said, by the formula (15.3).

Since the absolute value of $f(x)$ is expressed by the formula

$$(15.6) \quad |f(x)| = f_1(x) + f_2(x),$$

we have

$$(15.7) \quad \int_E |f(x)| dx = \int_E f_1(x)dx + \int_E f_2(x)dx.$$

Thus, the summability of $f(x)$ implies that of $|f(x)|$. Conversely, since $f_1(x) \leq |f(x)|$, $f_2(x) \leq |f(x)|$, the summability of $|f(x)|$ implies that of $f_1(x)$ and $f_2(x)$, therefore of $f(x)$. In short, the summability of a function $f(x)$ is logically equivalent to the summability of its absolute value.

It is often necessary to have the definition of Lebesgue

SURVEY OF THE THEORY OF INTEGRATION

integral over an unbounded set or interval. We define the Lebesgue integral of $f(x)$ over an interval $a \infty$ by the formula

$$(15.8) \quad \int_a^\infty f(x)dx = \lim_{b \rightarrow \infty} \int_a^b f_1(x)dx - \lim_{b \rightarrow \infty} \int_a^b f_2(x)dx,$$

provided that both limits exist—a necessary and sufficient condition for this is the absolute integrability of $f(x)$, i.e., the existence of

$$(15.8') \quad \int_a^\infty |f(x)| dx = \lim_{b \rightarrow \infty} \int_a^b |f(x)| dx.$$

If $f(x)$ is given on the entire real axis $(-\infty, +\infty)$, we define

$$(15.9) \quad \int_{-\infty}^{+\infty} f(x)dx = \int_{-\infty}^0 f(x)dx + \int_0^{+\infty} f(x)dx.$$

16. *Properties of the Lebesgue integral*

The Lebesgue integral is more general than the Riemann integral. Every function integrable R is also Lebesgue integrable, and its Lebesgue integral is equal to its Riemann integral. But a function not integrable R may be integrable in the sense of Lebesgue, e.g., the function $\varphi(x)$ (considered in §5) equal to 1 or to 0 according as x is rational or irrational. Since the measure of the rational points of the interval 01 is zero, while that of the irrational points of 01 is equal to unity, we have, for the Lebesgue integral,

$$(16.1) \quad \int_0^1 \varphi(x)dx = 0.$$

On the other hand, as we have seen in §5, the function $\varphi(x)$ is not integrable R .

The Lebesgue integral $\int_E f(x)dx$ is additive as to the set E of integration, whether E be expressed as the sum of a finite or an infinite number of measurable sets (com-

pletely additive function of sets). Specifically, if

$$E = E_1 + E_2 + \dots + E_m,$$

then

$$(16.2) \quad \int_E f(x)dx = \int_{E_1} f(x)dx + \int_{E_2} f(x)dx + \dots + \int_{E_m} f(x)dx;$$

also, if

$$E = E_1 + E_2 + \dots + E_n + \dots,$$

then

$$(16.3) \quad \int_E f(x)dx = \int_{E_1} f(x)dx + \int_{E_2} f(x)dx + \dots + \int_{E_n} f(x)dx + \dots,$$

the right member being a convergent infinite series.

Further, the Lebesgue integral is additive as to the integrand $f(x)$, for its expression as a finite number of terms: if

$$f(x) = f_1(x) + f_2(x) + \dots + f_m(x),$$

and each term $f_1(x), \dots, f_n(x)$ is summable on the set E , then $f(x)$ is summable on the set E , and

$$(16.4) \quad \int_E f(x)dx = \int_E f_1(x)dx + \int_E f_2(x)dx + \dots + \int_E f_m(x)dx.$$

Finally, we consider the Lebesgue integration of an infinite series.

It is possible to pass to the limit under the sign of Lebesgue integration under much broader conditions than

for Riemann integration. We have seen that, if the Riemann integrable functions $f_n(x)$ tend *uniformly* to $f(x)$ as $n \rightarrow \infty$, then $f(x)$ is Riemann integrable and

$$(16.5) \quad \int_a^b f(x)dx = \lim_{n \rightarrow \infty} \int_a^b f_n(x)dx.$$

On the other hand, if $f_n(x)$ merely stays bounded in its approach to $f(x)$, then an example (§6, second paragraph) has shown that $f(x)$ is not necessarily Riemann integrable.

With Lebesgue integration, uniform boundedness of the approaching summable functions $f_n(x)$:

$$(16.6) \quad |f_n(x)| < M, \quad f_n(x) \rightarrow f(x) \text{ as } n \rightarrow \infty,$$

M being a constant independent of n and of x , insures that the limit function $f(x)$ is summable and that

$$(16.7) \quad \int_E f(x)dx = \lim_{n \rightarrow \infty} \int_E f_n(x)dx.$$

Indeed, the weaker condition:

$$(16.8) \quad |f_n(x)| < g(x)$$

where $g(x)$ is a summable function independent of n , is sufficient for the summability of $f(x)$ and the relation (16.7).

An alternative form of statement is that any convergent infinite series of summable functions,

$$(16.9) \quad f(x) = u_1(x) + u_2(x) + \dots + u_n(x) + \dots,$$

is summable term by term:

$$(16.10) \quad \int_E f(x)dx = \int_E u_1(x)dx + \int_E u_2(x)dx + \dots + \int_E u_n(x)dx + \dots,$$

provided that the partial sums

$$(16.11) \quad S_m(x) = u_1(x) + u_2(x) + \dots + u_m(x)$$

stay uniformly bounded:

$$(16.12) \quad |S_m(x)| < M \quad (\text{all } m \text{ and } x),$$

or are uniformly dominated by a summable function $g(x)$:

$$(16.13) \quad |S_m(x)| < g(x) \quad (\text{all } m \text{ and } x).$$

17. *Fourier series and transforms*¹⁰

The *Fourier coefficients* a_n, b_n of a function $f(x)$, given in the interval $(0, 2\pi)$, are defined by means of the integrals of the product of $f(x)$ with trigonometric functions, namely:

$$(17.1) \quad a_n = \frac{1}{\pi} \int_0^{2\pi} f(x) \cos nx \, dx,$$

$$b_n = \frac{1}{\pi} \int_0^{2\pi} f(x) \sin nx \, dx,$$

$$(n = 0, 1, 2, 3, \dots).$$

These are made the coefficients of an infinite series:

$$(17.2) \quad \frac{a_0}{2} + \sum_{n=1}^{\infty} (a_n \cos nx + b_n \sin nx),$$

called the *Fourier series* of the function $f(x)$. The fundamental question associated with this series is that of its convergence to $f(x)$, whether in the ordinary sense, or in some extended sense, such as "convergence in the mean, of a given order p ". The last means that

$$(17.3) \quad \lim_{m \rightarrow \infty} \int_0^{2\pi} |f(x) - S_m(x)|^p \, dx = 0 \quad (p > 0),$$

where $S_m(x)$ denotes the m^{th} partial sum of the Fourier series:

$$(17.4) \quad S_m(x) = \frac{a_0}{2} + \sum_{n=1}^m (a_n \cos nx + b_n \sin nx).$$

¹⁰ For a technical and detailed account, see Titchmarsh, chap. XIII.

The case $p = 2$ is of particular importance throughout the theory.

Since integration figures in all the preceding formulas, the theory of Fourier coefficients and Fourier series evidently depends on the type of integration considered. Thus we have Fourier-Riemann and Fourier-Lebesgue coefficients and series.

The theory takes its neatest form when the integration is according to Lebesgue. Typical is the case of the Riesz-Fisher theorem: *if the series*

$$(17.5) \quad \frac{a_0^2}{2} + \sum_{n=1}^{\infty} (a_n^2 + b_n^2)$$

is convergent, then a function $f(x)$ exists, summable together with its square, such that the Fourier coefficients of $f(x)$ are precisely a_n, b_n . Furthermore,

$$(17.6) \quad \frac{1}{\pi} \int_0^{2\pi} f(x)^2 dx = \frac{a_0^2}{2} + \sum_{n=1}^{\infty} (a_n^2 + b_n^2).$$

Here Lebesgue integration is really essential for the validity of the theorem; a similar formulation with Riemann integration is not true.

In the proof of the Riesz-Fisher theorem, the required function $f(x)$ is obtained as the limit in the mean of order 2 of the partial trigonometric sums defined by (17.4).

Important in many applications are the *Fourier transforms* of a function $f(x)$, supposed defined in the interval 0∞ and summable in this interval according to the definition (15.8) or the criterion (15.8'). The Fourier cosine transform of $f(x)$ is

$$(17.7) \quad g(x) = \sqrt{\frac{2}{\pi}} \int_0^{\infty} f(t) \cos xt \, dt,$$

while its Fourier sine transform is

$$(17.8) \quad h(x) = \sqrt{\frac{2}{\pi}} \int_0^{\infty} f(t) \sin xt \, dt.$$

An important theorem is that the relation between $f(x)$ and $g(x)$, as well as that between $f(x)$ and $h(x)$, is a reciprocal one; i.e., as consequence of (17.7), we have

$$(17.9) \quad f(x) = \sqrt{\frac{2}{\pi}} \int_0^{\infty} g(t) \cos xt \, dt,$$

and, as a consequence of (17.8),

$$(17.10) \quad f(x) = \sqrt{\frac{2}{\pi}} \int_0^{\infty} h(t) \sin xt \, dt.$$

In case $f(x)^2$ is summable in 0∞ , then for purposes of the general theory of Fourier transforms, the integrals over 0∞ which serve to define them are most conveniently interpreted as limits in the mean of order 2. First, we write

$$(17.11) \quad g_a(x) = \sqrt{\frac{2}{\pi}} \int_0^a f(t) \cos xt \, dt;$$

then a unique function $g(x)$ always exists such that

$$(17.12) \quad \lim_{a \rightarrow \infty} \int_0^{\infty} |g_a(x) - g(x)|^2 \, dx = 0.$$

This $g(x)$ is, by definition, the Fourier cosine transform of $f(x)$. A similar definition gives the Fourier sine transform, $h(x)$.

The same relations of reciprocity (17.7), (17.9) hold for Fourier transforms in this sense of convergence in the mean. Also, we have

$$(17.13) \quad \int_0^{\infty} f(x)^2 \, dx = \int_0^{\infty} g(x)^2 \, dx = \int_0^{\infty} h(x)^2 \, dx.$$

V. REMARKS ON THE DENJOY INTEGRAL

18. *Inverse relation between differentiation and integration*

The most fundamental and typical fact of the differential and integral calculus is the inverse nature of differentiation

and integration, as expressed by the formulas:

$$(18.1) \quad \frac{d}{dx} \int_a^x f(x) dx = f(x),$$

$$(18.2) \quad \int_a^b F'(x) dx = F(b) - F(a) \quad (F'(x) = \frac{d}{dx} F(x)).$$

It is natural to examine the precise meaning and scope of application of these formulas when the functions and integrals involved are taken in the general senses that we have considered.

Formula (18.1) applies to the Riemann integral, certainly wherever $f(x)$ is continuous, therefore everywhere except at most on a set of zero measure—we say “almost everywhere”.

The same formula is proved in the text-books to apply to the Lebesgue integral almost everywhere on the interval ab .

The consideration of formula (18.2) is more difficult. First, we have the fact, due to Volterra,¹¹ that a function $F(x)$ may have a derivative $F'(x)$ everywhere finite-valued, but $F'(x)$ need not be integrable R. Therefore formula (18.2) does not hold universally with the Riemann interpretation of \int .

Further, derived functions $F'(x)$, everywhere finite-valued, exist, such that $F'(x)$ is not integrable in the sense of Lebesgue. An example is the function:

$$(18.3) \quad F'(x) = \left(x^2 \sin \frac{1}{x^2} \right)' = 2x \sin \frac{1}{x^2} - \frac{2}{x} \cos \frac{1}{x^2}$$

$$(x \neq 0), \quad F'(0) = 0.$$

The absolute value $|F'(x)|$ becomes unbounded as $x \rightarrow 0$

¹¹ See E. Hobson, *The Theory of Functions of a Real Variable*, 2^d ed., Cambridge, 1921-26, vol.1, p.461.

in such a way that the Lebesgue integral $\int_0^1 |F'(x)| dx$

does not exist.¹² It follows that $F'(x)$ itself is not summable over 01 (see end §15). Incidentally, the improper Rie-

mann integral of $F'(x)$ over 01 , defined as $\lim_{\epsilon \rightarrow 0} \int_{\epsilon}^1 F'(x) dx$,

evidently exists, being equal to $\sin 1$.

Thus, neither the Riemann nor the Lebesgue process of integration is sufficiently powerful to convert every finite-valued derived function $F'(x)$ back to its primitive $F(x)$. A process of integration adequate to solve this problem in all cases was invented by A. Denjoy,¹³ and called by him *totalisation*; his integral is called the *total* of $f(x)$ over the interval ab or the set E . With integration understood in this sense, the formula (18.2) holds universally for every finite-valued function $F'(x)$ that is the derivative of a primitive $F(x)$.

The effective formation of the Denjoy integral of a given function $f(x)$ involves a certain transfinitely repeated process, for the details of which we refer to the cited texts.¹⁴

VI. MULTIPLE INTEGRATION

19. *Riemann multiple integral*

The preceding theories of Riemann and Lebesgue integration can be extended to functions of two or any finite number of variables: $f(x, y)$, $f(x_1, x_2, \dots, x_m)$.

¹² Titchmarsh, p.342.

¹³ Denjoy's first note appeared in the Paris *Comptes Rendus* in 1912, and his more developed papers elsewhere in 1915, 1916, 1917.

¹⁴ Lebesgue, chap.X; Saks, chap.VIII; Kestelman, chap.9, §2; Denjoy part 2 chap.I.

SURVEY OF THE THEORY OF INTEGRATION

To fix the ideas, let $f(x, y)$ be any bounded function of the two independent variables x, y , defined in the rectangle

$$R: \quad a \leq x \leq b, \quad c \leq y \leq d.$$

Let this rectangle be subdivided by the lines

$$\begin{aligned} x &= a(\text{or } x_0), x_1, x_2, \dots, x_{m-1}, x_m(\text{or } b) \\ y &= c(\text{or } y_0), y_1, y_2, \dots, y_{n-1}, y_n(\text{or } d). \end{aligned}$$

Let r_{jk} denote the rectangle

$$x_{j-1} \leq x \leq x_j, \quad y_{k-1} \leq y \leq y_k,$$

and also the area of this rectangle:

$$(19.1) \quad r_{jk} = (x_j - x_{j-1})(y_k - y_{k-1}).$$

Let M_{jk}, m_{jk} denote respectively the least upper and greatest lower bounds of $f(x, y)$ in r_{jk} , and form the "upper and lower sums"

$$(19.2) \quad S = \sum_{j=1}^m \sum_{k=1}^n M_{jk} r_{jk}, \quad s = \sum_{j=1}^m \sum_{k=1}^n m_{jk} r_{jk}.$$

The lower bound of S for all possible partitions of R into rectangles r_{jk} is, by definition, the "upper Darboux integral":

$$(19.3) \quad \overline{\iint}_R f(x, y) dx dy = \underline{B} S;$$

it is also the limit of S as the "norm of the partition".

$$(19.4) \quad \Delta = \max_{j, k} (x_j - x_{j-1}, y_k - y_{k-1}),$$

tends to zero:

$$(19.5) \quad \overline{\iint}_R f(x, y) dx dy = \lim_{\Delta \rightarrow 0} S.$$

Similarly, the defining formulas of the "lower Darboux integral" are

$$(19.6) \quad \underline{\int \int}_R f(x, y) dx dy = \underline{B} s = \lim_{\Delta \rightarrow 0} s.$$

It can be shown that, as in the case of a single independent variable, we have always

$$(19.7) \quad \overline{\int \int}_R f(x, y) dx dy \geq \underline{\int \int}_R f(x, y) dx dy.$$

In the case of equality of upper and lower Darboux integrals, the function $f(x, y)$ is said to be integrable in the sense of Riemann over the rectangle R , and we write

$$(19.8) \quad \int \int_R f(x, y) dx dy = \underline{B} S = \overline{B} s = \lim_{\Delta \rightarrow 0} S = \lim_{\Delta \rightarrow 0} s.$$

Every continuous function $f(x, y)$ is Riemann integrable. A necessary and sufficient condition for the Riemann integrability of $f(x, y)$ is that the set of its points of discontinuity shall have zero two-dimensional measure, this notion being defined in the next section.

If the region K of definition of $f(x, y)$ is a general plane region bounded by one or more ordinary closed curves, then we can enclose K in a rectangle R with sides parallel to the axes, and define $\varphi(x, y)$ in the rectangle R as follows:

$$(19.9) \quad \begin{aligned} \varphi(x, y) &= f(x, y) \text{ if the point } (x, y) \text{ is in } K, \\ \varphi(x, y) &= 0 \text{ if the point } (x, y) \text{ is in } R - K. \end{aligned}$$

The Darboux and Riemann integrals of $f(x, y)$ over K are then defined as equal to those of $\varphi(x, y)$ over R .

20. Lebesgue multiple integral

To define the Lebesgue integral for a function $f(x, y)$ defined on a set E of points of the (x, y) -plane, we need

the notion of the plane or two-dimensional measure of any set of points (x, y) .

Let E denote any such set, contained in a given rectangle $R = (A \leq x \leq B, C \leq y \leq D)$. A system of axis-parallel rectangles $J = \{R_n = (a_n \leq x \leq b_n, c_n \leq y \leq d_n)\}$, $(n = 1, 2, 3, \dots)$, finite or denumerably infinite in number, will be called a "covering-system of the set E " if every point (x, y) of E is contained as interior or boundary point in some rectangle $R_{n(x, y)}$. The total area of the rectangles of the system J is the sum of the finite or infinite series

$$(20.1) \quad A(J) = \sum_n (b_n - a_n)(d_n - c_n),$$

and the lower bound of $A(J)$ for all possible covering-systems of E is called the "outer measure" of E :

$$(20.2) \quad \overline{m}(E) = \underline{B} A(J) \quad [\text{all covering-systems } J \text{ of } E].$$

The "inner measure" of E is defined by the formula

$$(20.3) \quad \underline{m}(E) = R - \overline{m}(CE),$$

where $R = (B - A)(D - C)$ denotes the area of the basic rectangle in which E is contained, and CE denotes the complementary set of E with respect to this rectangle.

Always, $\underline{m}(E) \leq \overline{m}(E)$. In case $\underline{m}(E) = \overline{m}(E)$, we call the set E "measurable", with measure $m(E)$, where $m(E) = \underline{m}(E) = \overline{m}(E)$.

The notion of measure being thus defined for two-dimensional sets (and in a similar way for m -dimensional sets), the Lebesgue integral of a bounded or unbounded measurable function $f(x, y)$ is defined by formulas exactly similar to those given in §§14, 15 for the case of a function of a single independent variable.

21. *Content of a two-dimensional set*

If we restrict the number of rectangles in a covering-system j of a given set E to be *finite* ($n = 1, 2, \dots, m$), then define

$$(21.1) \quad A(j) = \sum_{n=1}^m R_n = \sum_{n=1}^m (b_n - a_n)(d_n - c_n),$$

and

$$(21.2) \quad \bar{c}(E) = \underline{B} A(j) \quad [\text{all finite covering-systems } j \text{ of } E], \\ \underline{c}(E) = R - \bar{c}(CE),$$

the numbers $\bar{c}(E)$, $\underline{c}(E)$ so arrived at are called respectively the “outer and inner content” of E . We always have $\underline{c}(E) \leq \bar{c}(E)$; in case $\underline{c}(E) = \bar{c}(E)$, we write $c(E) = \underline{c}(E) = \bar{c}(E)$, and call this common value the “content” of the set E . We may also observe that, as in the one-dimensional case, $\bar{c}(E) \geq \bar{m}(E) \geq \underline{m}(E) \geq \underline{c}(E)$.

There is an interesting relation of the Riemann and Lebesgue integrals of a function $f(x)$ of a single variable to the content and measure of two-dimensional sets. Let $f(x)$, defined in the interval ab , be represented graphically in the plane (x, y) by the set of points

$$E: \quad a \leq x \leq b, \quad 0 \leq y \leq f(x)$$

—we suppose $f(x) \geq 0$, to fix the ideas. Then the Riemann integral of $f(x)$ is equal to the *content* of the two-dimensional set E , while the Lebesgue integral of $f(x)$ is equal to the *measure* of this set:

$$(21.3) \quad \int_a^b f(x) dx = c(E), \quad \int_a^b f(x) dx = m(E).$$

SURVEY OF THE THEORY OF INTEGRATION

REFERENCES

- 1 H. Lebesgue, *Leçons sur l'Intégration et la Recherche des Fonctions Primitives*, Paris, 1928.
- 2 S. Saks, *Theory of the Integral*, Warsaw-Lwów, 1937.
- 3 H. Kestelman, *Modern Theories of Integration*, Oxford, 1937.
- 4 A. Denjoy, *Introduction à la Théorie des Fonctions des Variables Réelles*, Paris, 1937 (especially, part 2, chap. I).
- 5 Ch.-J. de la Vallée Poussin, *Cours d'Analyse*, 3^d ed., Louvain-Paris, 1914.
- 6 Ch.-J. de la Vallée Poussin, *Intégrales de Lebesgue, Fonctions d'Ensemble, Classes de Baire*, Paris, 1916.
- 7 E. Titchmarsh, *The Theory of Functions*, Oxford, 1932, chaps. X-XIII.

THE FOUR COLOR PROBLEM

By PHILIP FRANKLIN

THE FOUR COLOR PROBLEM

CONTENTS

	PAGE
1. Introduction	53
2. Euler's Theorem	57
3. Topographic Maps	58
4. Regular Maps	59
5. The Five Color Theorem	61
6. Configurations Reducible with Four Colors	62
7. Minimum Irreducible Maps	64
8. Special Coloration Theorems	65
9. Special Classes of Colorable Maps	67
10. Methods Involving Vertices and Edges	68
11. The Problem of Tait, and Petersen's Theorem	71
12. More Dimensions	76
13. Surfaces of Higher Genus	77
14. One-sided Surfaces	80
15. Empires	82
16. The Number of Colorations	83
17. Mutually Contiguous Countries	83
18. Conclusions	84

THE FOUR COLOR PROBLEM¹

1. INTRODUCTION. Geographical maps of an area into political subdivisions, as of a continent into countries, or a state into counties, are usually colored. While we may use the same color for several distinct subdivisions, as we wish to emphasize the boundary lines we generally choose the colors in such a way that any two subdivisions which touch along a line have different colors. It was known to geographers for a long time that the maps met with in practice could be colored in this way without using more than four distinct colors, and that for some maps no smaller number would suffice.

The corresponding mathematical question had been formulated as early as 1850 by De Morgan, but was first put before a wide mathematical public in 1878, when Cayley² proposed it to the London Mathematical Society, in an address published in the Proceedings of that Society.

In the next year, in Volume II of the *American Journal of Mathematics*, a "solution" was published by A.B.Kempe.³ However, the problem was definitely unsolved again in 1890, when P.J.Heawood⁴ pointed out an error in Kempe's reasoning. He did show, by a revision of Kempe's proof, that five colors are always sufficient. This is the only result related to the four color problem applying to all maps with-

¹ Revised version of a lecture at the Galois Inst. of Math. at Long Island University.

² A.Cayley, "On the Colouring of Maps," *Proceedings of the London Mathematical Society*, v.9, 1878, p.148.

³ A.B.Kempe, *Am. Journal of Math.*, v.2, 1879, p.193.

⁴ P.J.Heawood, *Quarterly Journal of Math.*, v.24, 1890, p.332.

out restriction that has been proved. Thus, for ordinary maps, we know that five colors suffice. It is easy to show that four colors are sometimes necessary, but whether five are ever needed is an open question. The four color problem is that of showing that four colors suffice. It is, today, perhaps the simplest unsolved problem of mathematics. It is much simpler than the many famous unproved conjectures of number theory, since these involve such concepts as prime number or powers of an integer and at least a fair grasp of arithmetic. On the other hand, the map-coloring problem is intelligible to any kindergarten child equipped with an outline map and a box of crayons.

Mathematically considered, our maps may be drawn on a plane or sphere. Since the coloring problem is not affected by a deformation of the map, it belongs to the branch of mathematics known as topology,⁵ which is concerned with those geometrical properties of figures unchanged by deformations. While this field presents many difficult problems, the difficulties are usually due to such complications as two curves intersecting in an infinite number of points, related to the generality of the notion of continuous curve, or continuity. That part of topology in which continuity questions play a minor rôle is called combinatorial topology,⁶ and the coloring problem belongs in this restricted field. In fact, the boundary lines of the map may be deformed into straight lines or circular arcs.

Thus, for our purposes, an ordinary map may be replaced by a subdivision of a sphere into regions by a finite number of circular arcs such that no region touches itself along an

⁵ For general topological concepts, the article, "What Is Topology?" *Philosophy of Science*, v.2, 1935, p.39, as well as the references therein, may be consulted.

⁶ See, for example, M. A. St. Lagüe, "Les Réseaux," *Mémoires des. Sc. math.*, XVIII, Paris, 1926.

THE FOUR COLOR PROBLEM

arc. The arcs are the *sides* of the regions, and their end points the *vertices*. Two regions of like color may have one or more vertices in common, but they may not have a common side.

One might expect that, since the coloring problem has not been solved for ordinary maps, it would become even more difficult if we considered maps on surfaces more complicated than the sphere, as an anchor ring (surface of the rings of Saturn!) or the corresponding surfaces with more than one hole. Curiously enough, however, the coloring problem has been solved in all such cases which have been seriously attacked. Naturally, the methods fail in the plane case. The generalizations to three dimensions all prove to be trivial, since in higher dimensions there is no limit to the number of colors required.

Many European mathematicians have worked at the problem since Kempe and Heawood, two of the most recent being Errera⁷ and St. Lagüe.⁸ Of the many American mathematicians we may mention G.D.Birkhoff,⁹ Veblen,¹⁰ Brahana,¹¹ Ballantine, myself, Reynolds,¹² Frink,¹³ Whitney,¹⁴ Gill, Kagno,¹⁵ and Kittell. The endurance record of

7 A.Errera, "Exposé Historique du Problème des Quatre Couleurs," *Periodico di Matematiche*, v.7, 1927, p.20.

8 M. A. St. Lagüe, "Géometrie de Situation et jeux," *Mémorial des. Sc. math.*, XLI, Paris, 1929.

9 G.D.Birkhoff, "The Reducibility of Maps," *Am. Journal of Math.*, v.35, 1913, p.115.

10 O.Veblen, *Annals of Math.*, v.14, 1912-1913, p.86.

11 H.R.Brahana, "The Four Color Problem," *Am. Math. Monthly*, 1923, p.234.

12 C.N.Reynolds, Jr., "On the Problem of Coloring Maps in Four Colors," I, p.1 and II, p.427, *Annals of Math.*, v.28, 1927.

13 O.Frink, *Annals of Math.*, v.27, 1926, p.491.

14 Hassler Whitney, "A Theorem on Graphs," *Annals of Math.*, v.32, 1931, p.378.

15 I.Kagno, "Note on the Heawood Color Formula," *Journal of Math. and Physics*, v.14, 1935, p.228.

interest in the problem undoubtedly goes to P.J.Heawood, who became interested in the problem in the eighties while a student under Cayley, published his first paper¹⁶ on the problem in 1890, and several others including one¹⁷ as recent as 1935.

Using methods due to Kempe¹⁸ and Birkhoff¹⁹ to extend their results, in 1922 I succeeded²⁰ in showing that all maps of 25 or fewer regions can be colored in four colors, and in 1926 Reynolds²¹ improved the 25 to 27. Recently I found some further results which, combined with some reductions due to Errera²² and Winn,^{23,24} enabled me²⁵ to increase this number to 31. Winn²⁶ later extended this to 35.

What I propose to do here is to prove the result concerning five colors and illustrate how the problem may be simplified to the consideration of regular maps. I shall then define the notion of reducible configuration, establish some of the simplest of these, and describe several others. I shall mention certain alternative problems which are either

16 P.J.Heawood, *Quarterly Journal of Math.*, v.24, 1890, p.332.

17 P.J.Heawood, *Proc. London Math. Society*, (2) v.40, 1935, p.189.

18 A.B.Kempe, *Am. Journal of Math.*, v.2, 1879, p.193.

19 G.D.Birkhoff, "The Reducibility of Maps," *Am. Journal of Math.*, v.35, 1913, p.115.

20 "The Four Color Problem," *Am. Journal of Math.*, v.44, 1922, p.225.

21 C.N.Reynolds, Jr., "On the Problem of Coloring Maps in Four Colors," I, p.1 and II, p.427, *Annals of Math.*, v.28, 1927.

22 A.Errera, "Une Contribution au Problème des Quatre Couleurs," *Bull. de la Soc. Math. de France*, v.53, 1925, p.42.

23 C.E.Winn, "A Case of Coloration in the Four Color Problem," *Am. Journal of Math.*, v.49, 1937, p.515.

24 C.E.Winn, "On Certain Reductions in the Four Color Problem," *Journal of Math. and Physics*, v.16, 1938, p.159.

25 "Note on the Four Color Problem," *Journal of Math. and Physics*, v.16, 1938, p.172.

26 C.E.Winn, "On the Minimum Number of Polygons in an Irreducible Map," *Am. Journal of Math.*, v.62, 1940, p.406.

THE FOUR COLOR PROBLEM

equivalent to the four color problem, or at least would imply that four colors would be sufficient. Finally, I shall indicate some of the specialized theorems on coloration, and some of the generalizations. These will include fairly general theorems on 2, 3, 5, 6, and 7 colors, with only that for 4 lacking!

2. EULER'S THEOREM. For any map drawn on a sphere and all of whose regions are simply connected there is a relation between V the number of vertices, E the number of edges or sides, and F the number of faces or regions. By simply connected regions we mean those deformable into circles. The simplest example of a region which is not simply connected, but multiply connected, is a region deformable into the ring bounded by two concentric circles.

The relation in question is

$$V - E + F = 2. \quad (1)$$

It is universally known as Euler's theorem, because it was stated by Euler in 1752. However, it had been used by Descartes as early as 1640. It may be established by mathematical induction, noting that any map of the type in question may be constructed by starting with a sphere divided into two regions by two lines joining two points, and adding successively new points and lines. Each vertex added on a side already present increases V and E each by one, while each side added joining two vertices already present increases E and F by one each. Thus these operations do not change $V - E + F$. But the equation (1) holds for the original map, since $2 - 2 + 2 = 2$, so that it holds for any map derived from it.

The Euler result is frequently stated for a polyhedron, with the unnecessary restriction that it be convex, and the essential requirement that the faces be simply connected

tacitly assumed. It is true for any polyhedron with simple faces which may be deformed into a map on a sphere.

3. TOPOGRAPHIC MAPS. While we have introduced the Euler relation because it is useful, at some stage or other, in every attack on the four color problem or its generalizations, we will digress here for an application to topographic maps.²⁷ Suppose that a family of level lines is drawn on a sphere, as well as their orthogonal trajectories, or lines of greatest slope. For most of the sphere these will form small rectangles. Exceptions will occur at the tops of mountains, or points of maximum height, M , at the depths of depressions, or points of minimum height, D , and at the saddle points, S . For simplicity, let there be only a finite number of points M , D , and S , and let all the points S be simple saddle points. Let us consider each edge as contributing $1/2$ to each region it touches, and each vertex as contributing $1/4$ to the region having a square corner at it. Then for a rectangle, the contribution to $V - E + F$ is:

$$4(1/4) - 4(1/2) + 1 = 0. \quad (2)$$

Near a point M , or D , by omitting all inside a closed level line, as in Fig. 1a, we find for $V - E + F$:

$$n(1/2) - n(1/2) + 1 = 1. \quad (3)$$

Near a simple saddle point, we obtain essentially the configuration of Fig. 1b, and the contribution to $V - E + F$ is:

$$n(1/2) + 8(1/4) - (n + 8)(1/2) + 1 = -1. \quad (4)$$

²⁷ Compare A. Cayley, "Sur les courbes de niveau et les lignes de pente," *Phil. Mag.* (4), v.18, 1859, p.264, and also the author's paper, "Regions of Positive and Negative Curvature on Closed Surfaces," *Journal of Math. and Physics*, v.13, 1934, p.253, in particular pages 256 and 257.

THE FOUR COLOR PROBLEM

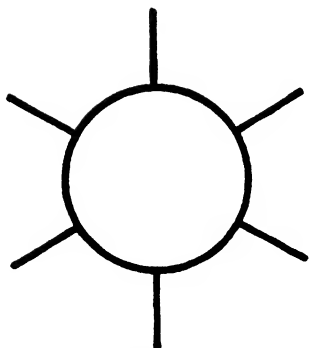


Fig. 1a

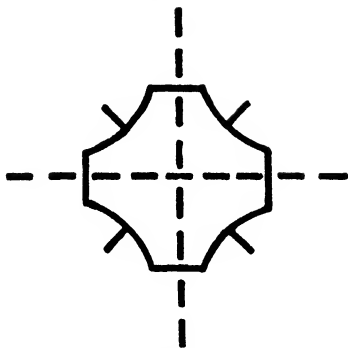


Fig. 1b

It follows from these results and the Euler relation that:

$$M + D - S = 2. \quad (5)$$

This shows that there are always at least two mountains or depressions, with additional ones if there are saddle points. An example is the sphere with saddle points at each pole, and four meridians as level lines. Here $S = 2$, $M = 2$, and $D = 2$.

4. **REGULAR MAPS.** The vertex of a map is called a triple vertex if three and only three edges meet there. A map all of whose vertices are triple, and all of whose regions are simply connected, is called a *regular map*. We shall now show that, in any coloring problem we may confine our attention to regular maps, since if all regular maps can be colored in N colors, then all maps can be colored in N colors.

First suppose that the map has some non-triple vertices. Then, if we cover each vertex at which n edges meet by a small circle we shall have a configuration similar to Fig. 1a. Then if we regard these small circles, or patches as Kempe called them, as new regions, we shall have a new map with n triple vertices in place of each n -tuple vertex. After

the new map is colored, we have only to remove the covering circles to obtain a coloration of the original map.

Next consider a map with one or more multiple regions. Consider, for example, a ring-shaped region, R , separating an inside map, I , from an outside map, O . If we could color I' , consisting of I surrounded by a single region R_1 , and O' , consisting of O surrounding a single region R_2 , we could color the original map. For, by a suitable interchange of two colors in one of the maps, say O' , we could arrange matters so that R_2 in O' had the same color as R_1 in I' . We need then merely color I as in I' , O as in O' and give R the common color of R_1 and R_2 to color the original map. This process, combined with a mathematical induction on the number of multiply connected regions, enables us to completely dispose of them.

We may deduce an important property of regular maps on a sphere from Euler's theorem. To do this, let A_n be the number of regions in the regular map with n sides. Then, obviously

$$\Sigma A_n = F. \quad (6)$$

We next count ends of edges by edges, vertices, and regions, and so find:

$$2E = 3V = \Sigma nA_n. \quad (7)$$

If we use the last two equations to eliminate V , E , and F from the equation (1), we obtain

$$(\Sigma nA_n)/3 - (\Sigma nA_n)/2 + \Sigma A_n = 2, \quad (8)$$

or

$$\Sigma(6 - n)A_n = 12. \quad (9)$$

Since the right member of this equation is positive, and

THE FOUR COLOR PROBLEM

the only positive terms in the left member result from values of n less than 6, it follows that *every regular map on a sphere contains at least one region of less than 6 sides.*

5. THE FIVE COLOR THEOREM. In order to prove the sufficiency of five colors, we introduce the notion of *reducibility*. We call a map *reducible*, if its coloration may be made to depend on the coloration, in the same number of colors, of a regular map with fewer regions. Any type of region, or combination of regions, is called a *reducible configuration* if its presence in a map renders it reducible. When we have five colors at our disposal, regions of two,

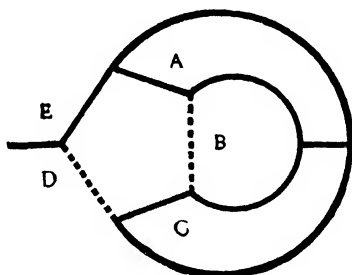


Fig. 2

three, or four sides are immediately reducible. For if we shrink such a region to a point, we obtain a map of fewer regions. If this last is colorable, so is the original map; since, when we restore the region, it will only be in contact with at most four distinct colors, so that there will be a fifth available for it.

But we may also reduce a pentagon when we have five colors. For, if the regions adjacent to a pentagon are in order a, b, c, d, e , either a and c do not touch each other, or if they do, as in Fig. 2, they prevent b from touching d . Thus, if an appropriate pair of edges of the pentagon are erased, a new map of fewer regions will result which has

no region touching itself along an edge. In Fig. 2 the dotted lines may be taken as the pair to be erased. If the reduced map is colored, and the pentagon then restored, the original map may be colored. For our construction has guaranteed that some pair of regions adjacent to the pentagon have the same color. Consequently, the pentagon touches at most four distinct colors, and a fifth is available for it.

Since we showed in Section 4 that every regular map must contain at least one region of five or fewer sides, it follows that every such map is reducible in five colors. Suppose now that there were a map not colorable in five colors.

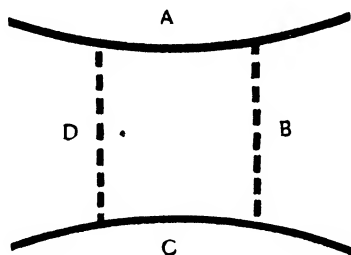


Fig. 3

Then there would be some regular map, M_k , with a minimum number of regions, say k . Hence every regular map with $k - 1$ regions would be colorable. But our reduction process enables us to color M_k by reducing its coloration to a map of $k - 1$ regions. This contradicts the assumption that M_k existed and proves the theorem that:

Every map on a sphere can be colored in five colors.

6. CONFIGURATIONS REDUCIBLE WITH FOUR COLORS. When we are limited to four colors, we may again consider a *reducible configuration* as any figure whose presence in a map renders it reducible. The argument of the preceding section shows that regions of two or three sides are reducible.

ble configurations, but breaks down for those of four sides. However, with four colors, a region of four sides is a reducible configuration.

We show this by arguing as follows. If the regions adjacent to the quadrilateral are in order a, b, c, d , either a and c do not touch one another, or if they do, they prevent b from touching d . Thus, if an appropriate pair of opposite edges of the quadrilateral are erased, a new map of fewer regions will result which has no region touching itself along an edge. In Fig. 3 the lines to be erased are dotted. If the reduced map is colored, and the quadrilateral then restored, the original map may be colored. For our construction has guaranteed that some pair of regions adjacent to the quadrilateral have the same color. Consequently, the quadrilateral touches at most three distinct colors, and a fourth is available for it.

Other reducible configurations are established by somewhat similar reasoning, except that there are more possibilities to be considered, and it is usually necessary to resort to the notion of *chains*^{28, 29} in two colors, and to consider the effect on a coloration of leaving one pair of colors fixed, and permuting the other pair in a connected portion of a partially colored map.

Among the simpler reducible configurations are non-triple vertices, multiply connected regions, and rings of four or fewer regions (Birkhoff³⁰). Some other known reducible configurations are:

(a) A ring of five regions not surrounding a single pentagon (Birkhoff³¹). This is fundamental in the proof of most

²⁸ A.B.Kempe, *Am. Journal of Math.*, v.2, 1879, p.193.

²⁹ G.D.Birkhoff, "The Reducibility of Maps," *Am. Journal of Math.*, v.35, 1913, p.115.

³⁰ *Ibid.*

³¹ *Ibid.*

of the other reductions.

(b) A pentagon adjacent to three consecutive pentagons (Birkhoff³²).

(c) A hexagon adjacent to three consecutive pentagons (the writer³³).

(d) A heptagon adjacent to four consecutive pentagons (Winn³⁴).

(e) A ring formed of an even number of pentagons and one or two adjacent additional regions (Winn³⁵).

(f) A ring formed of an even number of hexagons and zero or more pairs of adjacent pentagons surrounding one, two, or three contiguous regions (Errera³⁶).

The reductions (e) and (f) originated with Birkhoff³⁷ for rings surrounding a single region, without the additional regions or pairs of pentagons. The extension to them was due to the writer.³⁸ Later Errera³⁹ gave the extension for both types indicated in (f), and with certain restrictions, to any such rings. Finally Winn⁴⁰ removed the restrictions for (e).

7. MINIMUM IRREDUCIBLE MAPS. Most of the positive results for the four color problem so far have been obtained

32 G.D.Birkhoff, "The Reducibility of Maps," *Am. Journal of Math.*, v.35, 1913, p.115.

33 "The Four Color Problem," *Am. Journal of Math.*, v.44, 1922, p.225.

34 C.E.Winn, "On Certain Reductions in the Four Color Problem," *Journal of Math. and Physics*, v.16, 1938, p.159.

35 *Ibid.*

36 A.Errera, "Une Contribution au Problème des Quatre Couleurs," *Bull. de la Soc. Math. de France*, v.53, 1925, p.42.

37 G.D.Birkhoff, "The Reducibility of Maps," *Am. Journal of Math.*, v.35, 1913, p.115.

38 "The Four Color Problem," *Am. Journal of Math.*, v.44, 1922, p.225.

39 A.Errera, "Une Contribution au Problème des Quatre Couleurs," *Bull. de la Soc. Math. de France*, v.53, 1925, p.42.

40 C.E.Winn, "On Certain Reductions in the Four Color Problem," *Journal of Math. and Physics*, v.16, 1938, p.159.

THE FOUR COLOR PROBLEM

by using reducible configurations. The method of attack has been an attempt at a mathematical induction. That is, we assume the result false. This implies that there are some maps not colorable with four colors. In particular there would be one, at least, with a minimum number of regions. We call such a map a *minimum irreducible map*. If, then, any configurations are known to be reducible in four colors, such configurations must be absent from this minimum irreducible map. For, otherwise, this map would be colorable since its coloration would be reducible to that of a map of fewer regions, and all maps of fewer regions are colorable.

From the results stated at the beginning of Section 6, if it exists, the minimum irreducible map must be regular and have no region of less than five sides. Thus the only positive term in the left member of equation (9) is A_5 , so that the map contains at least 12 pentagons, and more if there are regions of more than six sides. It is by a consideration of the number of regions of more than six sides necessary to keep the pentagons and hexagons from forming known reducible configurations, including the additional pentagons necessitated by the presence of such regions, that the results of Reynolds, Winn, and myself, mentioned in Section 1 have been obtained. We recall that the best result is due to Winn⁴¹ and states that:

The minimum irreducible map contains at least 36 regions.

8. SPECIAL COLORATION THEOREMS. Several interesting results concerning special colorations have been proved. One is a *two color* theorem:

A necessary and sufficient condition for a map to be colored in two colors is that all the vertices be even.

⁴¹ C.E.Winn, "On the Minimum Number of Polygons in an Irreducible Map," *Am. Journal of Math.*, v.62, 1940, p.406.

This is related to the result, illustrated in certain Moorish Arabesques, that if a number of closed curves are drawn in a plane, the intersections may be so marked that each curve alternately passes over and under the curves it meets. The final pattern resembles a woven material, as seen in Fig. 4.

A familiar example of the two color theorem itself is a checker-board.

If, in a map drawn on a sphere, we select a point interior to each region, the capital of the country, and for each edge of the map on which two regions abut, draw a new line crossing it joining the corresponding capitals, we shall form a new map. The new map is called the *dual* of the first. The construction for a part of a map is shown in Fig. 5.

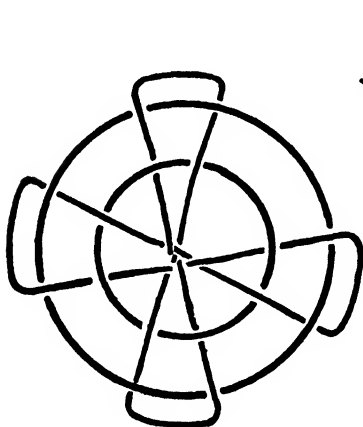


Fig. 4

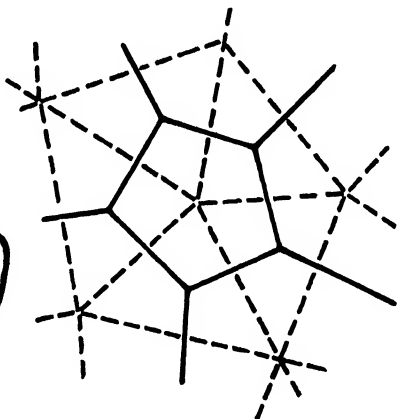


Fig. 5

The dual relation is reciprocal, if we regard the vertices of the first map as the capitals for the second. In number and association, vertices of one map correspond to regions for the dual map.

Suppose we start with a regular map, which is colored

THE FOUR COLOR PROBLEM

in three colors. We designate the colors by A , B , and C . Let us next put $+1$ at each vertex at which, following around from A in clockwise order we encounter ABC , and a -1 at each vertex at which, following around from A in clockwise order, we encounter ACB . Then by associating 1 and -1 with two colors given to the countries of the dual map in accordance with the marks at their capitals, we shall obtain a coloration of the dual map in two colors. This process may be reversed to give a coloration in three colors of the original regular map if the dual map is colored in two colors. Hence, as a consequence of the two color theorem, we have the *three color theorem*:

A necessary and sufficient condition for a regular map to be colored in at most three colors is that each region has an even number of sides.

A simple example is a map whose interior regions consist of hexagons arranged like the cells of a beehive. The map may be closed with one region of a large number of sides, and a number of quadrilaterals. If we regard the bricks as regions, a rectangular brick wall, as ordinarily constructed, provides a map topologically equivalent to that just described.

9. SPECIAL CLASSES OF COLORABLE MAPS. The theorems of the preceding section describe two large classes of maps which can be colored in not more than four colors. An additional class of this kind is given by the theorem:

A regular map with each region of $3n$ sides can be colored in four colors.

An example is the icosahedron, with all its faces triangular.

These results were all given by Heawood.⁴²

A known result of less restricted character is:

42 P.J.Heawood, *Quarterly Journal of Math.*, v.24, 1890, p.332.

A regular map containing at most one region of more than 6 sides can be colored in four colors.

This was proved by Winn,⁴³ by an extension and improvement of an argument given by Errera⁴⁴ who showed that any regular map containing no region of more than 6 sides was reducible, but not necessarily colorable.

So far no other extensive classes of maps which are easy to characterize have been proved to be colorable in four colors.

10. METHODS INVOLVING VERTICES AND EDGES. If we draw a small circle about a vertex of a regular map, and regard it as an added region, we obtain a new regular map. We refer to this process as triangulating a vertex. Now suppose that, for a given regular map, we select certain vertices, V' , which we leave unchanged, and that we triangulate the remaining vertices, V'' . If the triangulated map is colorable, so is the original map, since we have merely to omit the triangles about vertices V'' , and leave the coloring of the other regions as they were.

In particular, if the map can be triangulated in such a way that the new map has all its regions of $3n$ sides, it will be colorable by the theorem of the last section. The added triangular regions have 3 sides. Any other region of the triangulated map will correspond to a region of the original map, and will have 1 vertex for each vertex V' , and 2 vertices for each vertex V'' in this latter region. Thus our problem is to so select the vertices V'' that, when we put a figure 1 at vertices V' , and a figure 2 at vertices V'' , the sum taken around any region will be divisible by 3. We may replace the 2 by -1 , since $2 \equiv -1 \pmod{3}$.

⁴³ C.E.Winn, "A Case of Coloration in the Four Color Problem," *Am. Journal of Math.*, v.49, 1937, p.515.

⁴⁴ A.Errera, "Une Contribution au Problème des Quatre Couleurs," *Bull. de la Soc. Math. de France*, v.53, 1925, p.42.

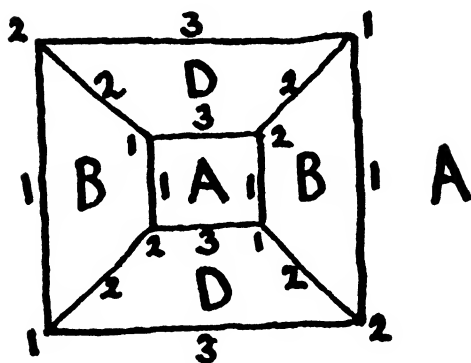


Fig. 6

In Fig. 6 we exhibit a simple map with vertices properly marked. This figure also illustrates a method of proceeding directly from the vertex marks to a coloration of the original map. We mark the edges with the figures 1, 2, or 3 (or 1, -1, 0) in such a way that all three figures occur at each vertex, the sense of rotation, 1, 2, 3 being clockwise at vertices marked with 1, and counter-clockwise at those marked with 2. The condition on the vertex marks that the sum around any region be divisible by 3 is equivalent to consistency in these edge marks. Thus we may mark any edge 1, and then working out from this mark the other edges, guided by the vertex marks and the rule just given.

Finally, we associate the three groups of pairs of colors and their complements with the numbers 1, 2, 3. That is, if we denote the four colors by A , B , C , and D , we consider 1 as AB or CD , 2 as AC or BD , and 3 as AD or BC . We now label any region as A , and use this table to color the adjacent regions, and then the ones adjacent to them, etc. Working out in this way, no inconsistency will arise if the edges are colored as described.

Conversely, the process just described may be reversed to obtain marks 1, 2, 3 for the edges, and hence marks 1, 2

for the vertices subject to the conditions previously mentioned. Thus our discussion proves:^{45,46}

The problem of coloring regular maps in four colors is equivalent to the problem of placing 1 or 2 (- 1) at each vertex in such a way that the sum, taken about each region, is divisible by 3.

This gives a theoretical method of coloring any given map. For the marks for the vertices may be taken as variables, and the conditions then lead to a system of linear equations. The complete solution of these equations may be found, and reduced modulo 3. Any solution with no variable zero gives a coloration, and if no such solution exists there is no coloration. This method is not practical even for simple maps. For example, the dodecahedron may be easily colored empirically but the present method leads to a system of 12 equations in 20 variables.

There are two difficulties met in the attempt to prove that suitable solutions always exist. One is the ruling out of zero values, which gives the problem a quadratic character. The other is the algebraic complication of the geometric condition that the equations should result from some map.

We may consider the problem of marking the edges with 1, 2, and 3 as that of coloring the edges in three colors, since the condition that all three appear at each vertex is equivalent to requiring that no two edges which meet at a vertex have the same color. Thus:

The coloring of the regions of a regular map in four colors is equivalent to the coloring of its edges in three colors.

In place of the edges, we may use an edge map, all of whose regions are four sided. Such a map may be obtained

45 P.J.Heawood, *Quarterly Journal of Math.*, v.24, 1890, p.332.

46 O.Veblen, *Annals of Math.*, v.14, 1912-1913, p.86.

THE FOUR COLOR PROBLEM

from a regular map by joining a point, or capital, of each region to the vertices of that region, and then erasing all the edges of the original map. We note that such edge maps are not regular. For this reason their coloration in three colors does not follow from the last theorem of section 8. Every map with all its regions four sided is not an edge map, and it is not easy to characterize those maps which are edge maps.

11. THE PROBLEM OF TAIT, AND PETERSEN'S THEOREM. As a preliminary step to coloring the edges in three colors, we may mark one set, say those with figure 1. The remaining edges will form one or more circuits. That a set of edges having just one end point of an edge at each vertex could be found in any regular map was proved by Petersen⁴⁷ in 1891. His proof has been simplified by Brahana⁴⁸ and Frink.⁴⁹ Petersen's theorem really applies to all third order linear graphs, satisfying a mild condition which in the regular map is a consequence of no region touching itself. A third order linear graph is the figure obtained by taking a finite number of points (not necessarily in a plane) and joining them by line segments in such a way that exactly three segments end at each point or vertex. The condition is that not more than three such segments form the only connection between two portions of the graph. Under these conditions, the theorem of Petersen asserts that it is always possible to color the segments in two colors, such that one end of the first color and two ends of the second color abut at each vertex.

The graph shown in Fig. 7 proves that the added condition is necessary. A graph to which the theorem applies is

47 J. Petersen, *Acta Mathematica*, v.15, 1891, p.193.

48 H.R. Brahana, *Annals of Math.*, v.18, 1917, p.59.

49 O. Frink, *Annals of Math.*, v.27, 1926, p.491.

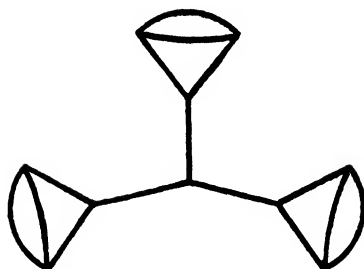


Fig. 7

shown in Fig. 8a. For example, we may take the two pentagons 1, 2, 3, 4, 5 and 6, 7, 8, 9, 10 as two circuits, passing twice through each vertex. This graph cannot be colored in three colors. It cannot be put in a plane, though it may be put on an anchor ring. One way of doing this is shown in Fig. 8b. The anchor ring is to be formed from the rectangle by first bringing the top around to touch the bottom, making a tube, and then bending the tube around to form the anchor ring. Thus the points on the opposite side of this rectangle should be regarded as the same. On the anchor ring, this graph makes a map with three pen-

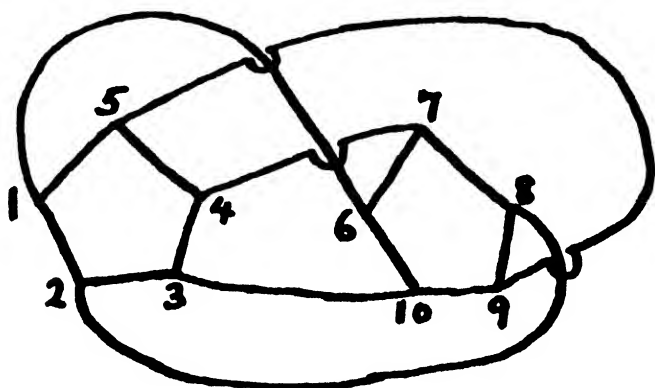


Fig. 8a

THE FOUR COLOR PROBLEM

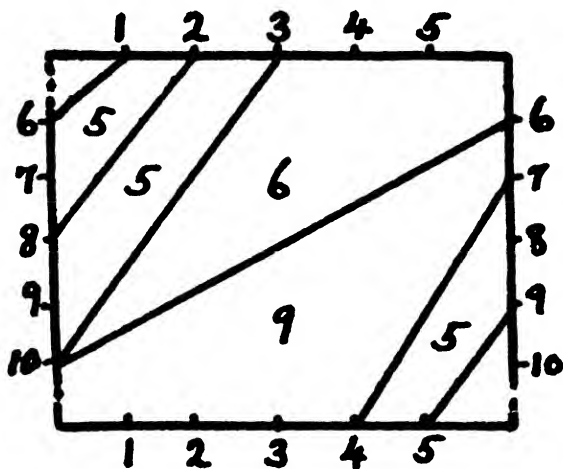


Fig. 8b

tagons, marked 5, one hexagon, marked 6, and one nonagon, marked 9 in Fig. 8b.

The example just given shows that to establish the edge coloration necessary for the four color problem, use must be made of the fact that the graph is plane. One simple geometric condition for a graph to be capable of representation in a plane was given by Kuratowski, and found inde-

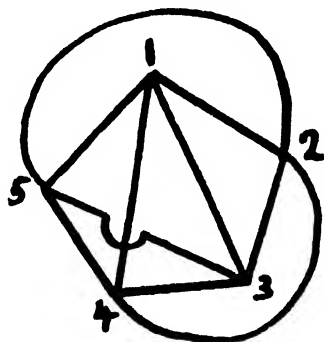
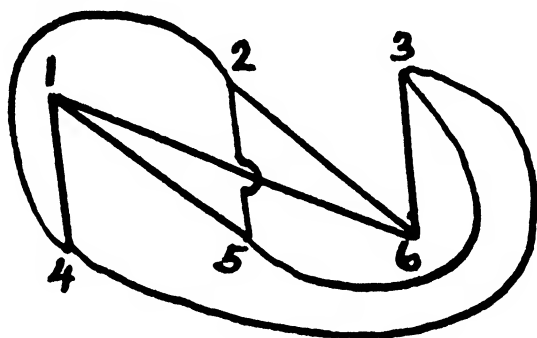


Fig. 9a

pendently by Frink. It is that the graph does not contain as a part of itself one of the configurations shown in Figs. 9*a* and 9*b*. The first consists of a pentagon with all its diagonals. The second consists of two triplets of points 1, 2, 3 and 4, 5, 6, each point in one triplet being joined to each point of the other. It is the basis of a familiar puzzle, to join each of three houses with each of three wells (or in a modern version to a gas, water, and electricity plant) without having any of the pipes cross. The puzzle is impossible on a plane, and only becomes possible if the pipes may be constructed on the surface of an anchor ring. The

Fig. 9*b*

configuration 9*b* is present in Fig. 8*a*, for example by taking 1, 7, 10 and 2, 5, 6 as the two triplets.

If, on decomposing a graph of the third order by Petersen's theorem, the circuits are all even, they may be broken up further to give a coloring of the edges of the graph in three colors. The difficulty comes when there is one or more pairs of odd circuits, like the pentagons of Fig. 8*a*.

Sir William Hamilton noticed that the edges of a dodecahedron could all be traversed by one circuit, and made a puzzle based on this fact. Tait conjectured that this

THE FOUR COLOR PROBLEM

held for any convex polyhedron with triple vertices. If this conjecture is correct, some condition implied by the convexity is necessary, as the example given in Fig. 10 shows. In this case there is essentially only one way of coloring the edges in three colors, namely as indicated. But, for this the 1-2, 2-3, and 3-1 circuits all break up into two parts. The edge coloration of Fig. 6 is also of this type as it stands. However, here this may be changed. For example, we may interchange 1 and 2 on one of the circuits

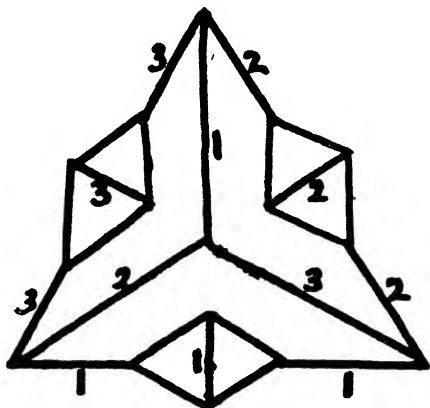


Fig. 10

and so obtain the edge coloration of Fig. 11, in which the 3-1 and 2-3 circuits each pass through all the vertices. Tait's conjecture would imply the four color theorem, but not conversely. If Tait's conjecture is false, there may be comparatively simple examples which disprove it without having rings of two regions as in Fig. 10. It would be interesting to find one.

A result in the nature of a dual to Tait's conjecture was established by Whitney,⁵⁰ who showed that in the dual
50 Hassler Whitney, "A Theorem on Graphs," *Annals of Math.*, v.32, 1931, p.378.

of a regular map, containing no two-sided regions, there must exist a circuit, composed of edges, which passes through all the vertices. For the original map this yields a circuit which passes once and only once through each capital. Whitney used this result to develop a canonical form for the dual of a regular map, which starts with a number of points on a circumference, one for each capital, and joins them by two series of diagonals, one inside the circle, and one outside the circle, to form triangular regions of the dual map.

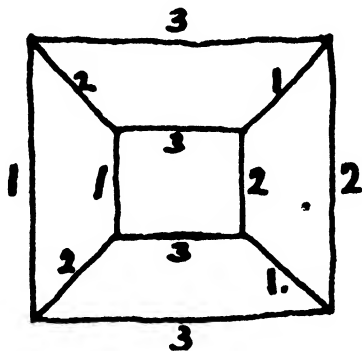


Fig. 11

12. MORE DIMENSIONS. If we increase the dimension of the color problem, and try to color volumes in space, the number of colors required may be unlimited. We should expect this, since any number of long strings, each of a different color, can be so braided that each touches every other. Another simple example is obtained by taking two layers of long square prisms, the top layer at right angles to the bottom layer, fastening the n th prism of the top layer to the n th prism of the bottom layer, and regarding each such pair fastened together as one block. Each block then touches every other. The number of colors remains un-

THE FOUR COLOR PROBLEM

limited, even if the regions be required to be convex.

13. SURFACES OF HIGHER GENUS. A better generalization is obtained by drawing our maps on a surface topologically different from the sphere. For such surfaces, the Euler relation (1) must be replaced by

$$V - E + F = K, \quad (10)$$

where K is the characteristic of the surface. For a surface obtained by removing p circular holes from a large disc, or fastening p handles on a large sphere, $K = 2 - 2p$. The number p is called the genus of the surface. We have $p = 0$ for the sphere, 1 for the anchor ring, and 2 for the sphere with two handles. The corresponding values of K are 2, 0, -2 . Thus K is less than or equal to zero for the surfaces now under discussion.

By the reasoning used to deduce (9) for the sphere, we find:

$$\Sigma(6 - n)A_n = 6K. \quad (11)$$

Heawood⁵¹ used this to determine the probable value of the number of colors needed for a map on any surface of higher genus. We observe that those arguments in sections 5 and 6 which depended on a closed circuit separating a map on a sphere into two parts are not applicable to the higher surfaces. However, it is still true that the color problem for arbitrary maps is reducible to that of coloring regular maps. Also the earlier argument shows that a map is reducible with $r - 1$ colors if it contains any region of $(r - 2)$ or fewer sides. Suppose, now, that a map on a surface with characteristic K is not reducible in $r - 1$ colors. Then $\Sigma A_n \geq r$, since the map must contain more than $r - 1$ regions. Also each n exceeds $r - 2$, and hence

51 P.J.Heawood, *Quarterly Journal of Math.*, v.29, 1898, p.270.

$\geq r - 1$. Thus $n - 6 \geq r - 7$. Let us assume that $r \geq 7$. Then $\Sigma(n - 6)A_n \geq (r - 7)r$. This may be combined with equation (11) to give:

$$-6K \geq (r - 7)r, \text{ or } r^2 - 7r + 6K \leq 0. \quad (12)$$

The equation $r^2 - 7r + 6K = 0$ has one root zero (for $K = 0$) or negative (for $K < 0$), and the other root positive. Let P be the greatest integer contained in this positive root, so that

$$P \leq \frac{7 + \sqrt{49 - 24K}}{2} < P + 1. \quad (13)$$

Then, for $r = P + 1$, $r^2 - 7r + 6K > 0$, which contradicts (12). This shows that no map is not reducible in $r - 1 = P$ colors. Consequently all maps are reducible, and hence colorable in P colors. Since (13) shows that $P \geq 7$ for $K \leq 0$, the assumption $r \geq 7$ is no restriction.

Let us next suppose that a map requires s colors because it contains, possibly with other regions, s regions each of

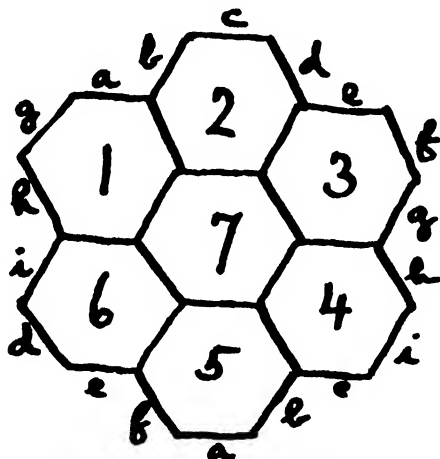


Fig. 12

THE FOUR COLOR PROBLEM

which touches all of the $s - 1$ other regions distinct from itself. Then each of these regions has $n \geq s - 1$, and since $\Sigma A_n \geq s$, we find

$$s^2 - 7s + 6K \leq 0, \quad (14)$$

by the reasoning used to deduce the inequality (12). Thus, the number P given by the relations (13) is the largest number for which there could be an example of this type with $s = P$.

If, for any value of K , we actually construct an example of this type, the color problem is completely solved for

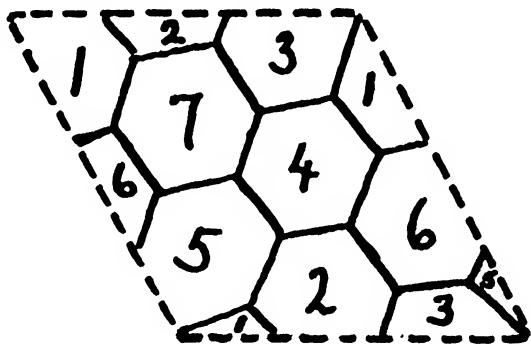


Fig. 13

that value of K . For, the example shows that P colors are sometimes necessary, and the earlier argument shows that P colors are always sufficient, so that P is the *chromatic* number. For $K = 0$, the anchor ring, the example requiring 7 colors was given by Heawood. A symmetric form of it is given in Fig. 12, which becomes an anchor ring if the sides with corresponding letters are brought into coincidence. This is more readily seen from Fig. 13, where the anchor ring is formed by combining opposite sides of the parallelogram, analogous to the process described in Section 11

for Fig. 9*b*. This proves that *the chromatic number for an anchor ring is seven*, a seven color theorem.

Examples for other cases were given by Heffter, so that for surfaces of genus 1, 2, 3, or 4 the chromatic number is known to be correctly given by the number P of the relation (13). If we denote the corresponding characteristic K by a subscript, we have:

$$P_0 = 7, P_{-2} = 8, P_{-4} = 9, P_{-6} = 10. \quad (15)$$

14. ONE-SIDED SURFACES. The surfaces of the last section have been two-sided, or orientable. The first term refers to the fact that, for example, a hollow sphere or anchor ring has an outer side and an inner side. The second term refers to the fact that a small clock face may be placed in the surface at a point in two ways, one clockwise for an observer on the outer side and one clockwise for an observer on the inner side, so that we may characterize a consistent sense of rotation for the surface as a whole. A one-sided surface, or non-orientable surface is one in which a small clock face may be so moved that it returns to the starting point with reversed sense. A band with a single twist in it is an illustration. Closed surfaces of this type in three dimensions necessarily have self-intersections, but they may be imaged on polygons analogous to the way we imaged an anchor ring on a rectangle. The simplest example is a projective plane, obtained by identifying opposite points of a circle, or polygon. This is the plane of projective geometry, and the change of orientation when we pass through infinity along a straight line explains why when we move from one end of an asymptote, with, say, a hyperbola on our left, we meet the hyperbola on our right at the other end. For the projective plane, $K = 1$ in the equation (10), and by fastening a number of knobs, each

THE FOUR COLOR PROBLEM

formed by cutting a single circular hole in a projective plane, to a sphere, we may form a surface of characteristic $K = 2 - k$, where k is the number of knobs. Thus while the characteristic is always even for a two-sided surface, it may have the value 1, 0, or any negative integer for a one-sided surface.

By reasoning as we did for the sphere, we may show that every map in the projective plane contains regions of five or fewer sides. Hence it is surely colorable with six colors. In Fig. 14 we give the example, due to Tietze, which shows

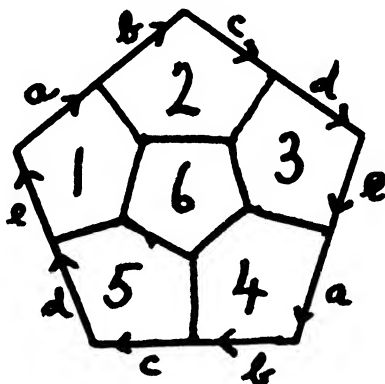


Fig. 14

that six colors are needed. Thus *the chromatic number for the projective plane is six*, a first six color theorem.

The Heawood argument given in section 13 may be applied to one-sided surfaces. It gives 7 for the one-sided surface with characteristic zero, sometimes called the Klein bottle. However, by a special argument I was able to show⁵² that only six colors are needed on this surface, and that *the chromatic number for the Klein bottle is six*. This is our second six color theorem.

⁵² "A Six Color Problem," *Journal of Math. and Physics*, v.13, 1934, p.363.

The examples which show that the Heawood number is correct have been given in several of the earlier cases by Kagno.¹⁵ Thus, using a prime to denote that the surface is one-sided, and as before denoting the characteristic K by a subscript, we have for the chromatic number:

$$P'_0 = P'_1 = 6, P'_{-1} = 7, P'_{-2} = 8, P'_{-4} = 9. \quad (16)$$

15. EMPIRES. Another way of generalizing the color problem is to consider a map drawn on a sphere, but composed of empires, with the further demand that a mother country and its colonies are similarly colored. If there is no limit to the number of colonies in each empire, the number of colors needed may be arbitrarily large. If no empire

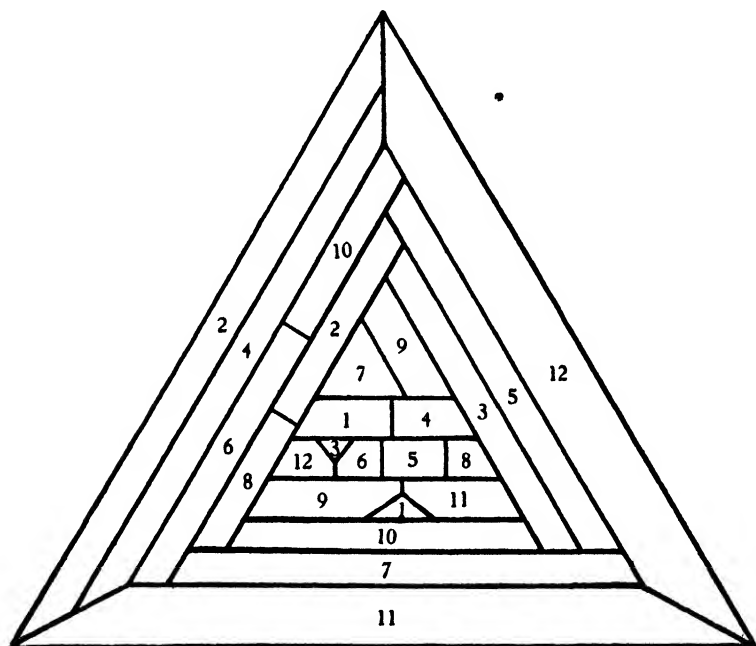


Fig. 15

THE FOUR COLOR PROBLEM

consists of more than r separated pieces, Heawood⁵³ proved that at most $6r$ colors are needed. He showed that for $r = 2$, the limit 12 is correct. One such map of twelve countries, each of which has one colony and one mother country is shown in Fig. 15. This cannot be colored in less than 12 colors, as shown, since each color touches every other one.

The formula is probably correct for greater values of r , but incorrect by 2, if the four color hypothesis is correct, and otherwise by 1, from the five color theorem, for the case $r = 1$.

Similar formulae were found by Heawood⁵⁴ for surfaces other than the sphere, which are probably correct though tedious to verify.

16. THE NUMBER OF COLORATIONS. A polynomial expression in p , which represents the number of ways of coloring a map in p colors, was given by Birkhoff.⁵⁵ In a number of later papers he and Whitney have studied the properties of this polynomial, with a view to showing that it does not vanish for $p = 4$. However, while the polynomial is quite tractable for other values of p , and its values have even been interpreted for fractional values, it has not yet been shown that the polynomial cannot have 4 as a root.

17. MUTUALLY CONTIGUOUS COUNTRIES. The color problem is frequently confused with the problem of finding the maximum number of regions, each of which touches all the remaining ones. This is related to the color problem, but only in the sense that if s regions of this type can be found, the chromatic number is at least s . For the sphere, the maximum value of s is known to be 4. It is an interesting

⁵³ P.J.Heawood, *Quarterly Journal of Math.*, v.29, 1898, p.270.

⁵⁴ *Ibid.*

⁵⁵ G.D.Birkhoff, "A Determinant Formula for the Number of Ways of Coloring a Map," *Annals of Math.*, v.14, 1912, p.42.

fact that for all the surfaces for which the color problem has been solved, the chromatic numbers of (15) and (16) are equal to the maximum number of mutually touching regions.

18. CONCLUSIONS. Since the four color problem has not been solved, in a mathematical sense it is not known whether it is true or not. However, practically all those who have worked on the problem are inclined to guess that it is true.

An admittedly crude probability argument was given by Heawood.⁵⁶ While the probability of the truth of a theorem, at a given stage of our knowledge, only has meaning if we consider it as one of a large class of theorems, the notion of the probability that a random map can be colored in four colors is a little easier to define. The argument indicates that the probability of being unable to color a regular map of F regions with four colors is of the order of $(1-3^{-F})^{2^F}$, or when F is fairly large $e^{-(\frac{1}{3})^F}$, approximately.

This is less than 1 in $10^{10,000}$ when F exceeds 35, so that if the probability argument is valid, and uncolorable maps exist, they are not easy to find.

Another question is whether the prospects seem bright for a solution of the problem, or at least a rapid advance in the number of regions known to be in an irreducible map. On this question we must be pessimistic, since the number has only risen to 35 in half a century. Also, at various stages, examples of maps of relatively few regions were given which, although colorable, escaped all the reductions so far known. These showed that further progress could only be made by finding new reductions.

⁵⁶ P.J.Heawood, *Proc. London Math. Society*, (2) v.40, 1935, p.189.

THE FOUR COLOR PROBLEM

Finally we note that, although the number 35 is small, it would not be easy to verify the result by experiment, since the number of topologically distinct maps on the sphere with 35 or fewer regions undoubtedly exceeds 10^{35} .

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

CHARLES SANDERS PEIRCE
AS A PIONEER

By CASSIUS JACKSON KETSER

CHARLES SANDERS PEIRCE AS A PIONEER

CONTENTS

Biographical Introduction	91
The Father of Pragmatism	92
Infinite	98
Propositional Function	102
Paradoxes	105
Relations	107

CHARLES SANDERS PEIRCE AS A PIONEER¹

BIOGRAPHICAL INTRODUCTION

THE aim of this talk is to tell you a little about some of the ideas that Charles Sanders Peirce gave or helped to give us and the world. Ideas, you know, are givers of light if the inner eye be fit, and, like stars, they differ in glory. It sometimes happens that folk are color-blind. I trust that none of us are idea-blind. Our being here today is some evidence that we are not.

Peirce was born in 1839. He died in 1914, just in time, you note, to be spared the pain of the World War and the dismal days that have followed.

In the years of his life there were many eminent men in the world, among them great men of science and great mathematicians. Perhaps it would interest you to list the names of a score or so of the most distinguished ones among them and at the same time to list their principal achievements. The task would not be a too easy one but it would be good fun and might add a bit to your wisdom. You might desire also to find out whether all of the great men in question were appreciated and duly honored in their day. Our pioneer was not, and he knew it. "I am a man," he once wrote, "of whom critics have never found anything good to say." But that is another matter. The critics are dead, most of them, and forgotten, while Peirce's fame is growing.

The father of our hero was Benjamin Peirce, long a dis-

¹ Lecture by C. J. Keyser at The Galois Institute of Mathematics, May 18, 1935.

tinguished member of the faculty of Harvard University. The son was graduated at Harvard, but he said of his father: "He educated me, and if I do anything it will be his work." The father was also a pioneer but his pioneering was in Mathematics while that of the son was in the general field of Logic. Benjamin Peirce's principal contribution to mathematics was his "Linear Associative Algebras" (*American Journal of Mathematics*, Vol. 4). Charles, too, was a mathematician but, as already said, his dominant interest was Logic.

But are not Mathematics and Logic the same thing? To say yes, as some have done, one must greatly broaden the current conception of mathematics or else greatly narrow that of logic, or do both. Now, to define two terms in such a way that they shall mean the same and then to assert solemnly that they do mean the same is not a very weighty contribution to knowledge! Charles Peirce was very far from regarding logic and mathematics as the same. Let us glance at the matter.

In the paper cited, Benjamin Peirce gave his now famous definition of mathematics as being "the science which draws necessary conclusions." Note the word "draws" and the word "necessary." The son approved of his father's definition, but he pointed out that to *draw necessary* conclusions is one thing, that to *draw* conclusions is another, and that the *science* of *drawing* conclusions is a third. And this science, said he, is Logic.

What *is* a necessary conclusion? Here we bump into the theory of Probability, which you may or may not have studied a little. What is the probability that a stated conclusion is warranted by the stated evidence? The probability may be zero or 1 or some fraction, any fraction, between. When and only when the probability is 1, the

conclusion is "necessary." At once you see that relatively few conclusions are "necessary." It is evident that to draw such conclusions is one thing; that to investigate the conditions for drawing them is a different *kind* of thing; and that to investigate the necessary and sufficient conditions for drawing conclusions, no matter what their probability, differs from actually drawing necessary conclusions both in kind and in scope. It plainly follows that, according to the conception of mathematics held by both the father and the son, and the conception of logic held by the son, mathematics and logic differ not only in kind but vastly in scope, and that logic's scope is the greater.

What has been said may enable you to glimpse why it was that Peirce regarded Logic as unsurpassed by any other subject in difficulty or dignity or range or human significance. He brought to the study of it a quite unsurpassed measure of the requisite type of genius, and devoted to it greater energy perhaps and more years than any other man in the annals of science. Peirce's logical researches inevitably led him to ponder and discuss almost every kind of fundamental question met with in the history of thought. In the course of his long life he published many papers, for the most part without material compensation. Of his unpublished papers there were at his death several hundred. Fortunately, an edition of his *Collected Papers*, constituting ten large volumes, is being put forth by the Harvard University Press, under the editorship of Dr. Charles Hartshorne and Dr. Paul Weiss. The first five volumes have already appeared. If you cannot afford to own the edition, you can at least request your branch of the Public Library to procure it for your use.

In the early years of The Johns Hopkins University Peirce held there, for a short time only, a lectureship in

GALOIS LECTURES

Logic. There was no separate department of logic and Peirce's lectures were given under the auspices of the department of mathematics, of which James Joseph Sylvester was then the head. The following story was told me the other day by a neighbor of mine, Mr.A, who, while pursuing mathematics under Sylvester's guidance in those days, had at the same time attended the course of Peirce's lectures in logic. One day Sylvester summoned Mr.A, for whom he had the fondness of a friend, and said to him: "You have listened to Mr.Peirce's lectures. Tell me about them. How have they impressed you?" Mr.A explained at some length that the lectures were always substantial, often very subtle, never trite, not easy to follow, frequently so lacking in clearness that the hearers were quite unable to understand, "but," added Mr.A, "there can be no question that Mr.Peirce is a man of genius." Thereupon Sylvester, who had been listening in silence,*said quite impulsively: "Well! If he is a genius, isn't that enough? Isn't it men of genius that we want here?"

Why Peirce did not continue at Hopkins I do not know. He never held a professorship there or at Harvard or elsewhere. Perhaps that fact was a blessing in disguise, but it was certainly attended by some serious disadvantages. "All my life," wrote Peirce, "my studies have been cruelly hampered by my inability to procure necessary books." How he contrived to do so vast an amount of exceedingly difficult work is something of a mystery.

He did most of it in the solitude of his home at Milford, Pennsylvania, where he lived many years, where he often had but scant means of subsistence and might have perished of starvation but for the intervention of some friends with sense enough to discern the quality of his genius, and where he died.

CHARLES SANDERS PEIRCE AS A PIONEER

Here was a man who immeasurably increased the intellectual wealth of our world and who nevertheless was in danger of starving in what was then the richest and vainest of lands. To think of it is enough to make the blood of any decent American boil with chagrin, indignation, and vicarious shame.

Peirce's widow died only a few months ago. In *Science* for November 16, 1934, Joseph Jastrow has written a beautiful and touching tribute to Mrs. Peirce's devotion and loyalty to the memory of the great man who was her husband. Do not fail to read it.

Let us now turn to some of the matters in which Peirce was a pioneer. It will be necessary to select, for there will not be time to mention all of them. Neither will it be possible to dwell at length on any of them.

The Father of Pragmatism

If your studies have made you fairly acquainted with modern philosophy, you can hardly have failed to encounter the term Pragmatism and you doubtless know something of the philosophy or the philosophic movement or the philosophic method which the term denotes. By its birth and main development Pragmatism is an American product, certainly not the weightiest but undoubtedly the most popular and appealing contribution of our land to philosophic thought. It sprang up here quite suddenly a little less than forty years ago, fairly flamed into existence, quickly over-ran the country, soon crossed the Atlantic and created quite a stir in philosophical circles both in Great Britain and on the continent of Europe. If you are not familiar with the pragmatic movement, the easiest and best way to gain a good knowledge of it is by reading William James's charming book on *Pragmatism*, for of all

the advocates, elaborators, and interpreters of the doctrine, James was the most engaging, illuminating, and powerful.

I trust it will be sufficient to remind you that Pragmatism is distinguished among rival philosophies by the answers it offers to three great questions, which are fundamental. These may perhaps be stated as follows:

- (1) What is the best (perhaps the only) way to ascertain the *meaning* of an idea or concept or proposition or doctrine?
- (2) What is the best (perhaps the only) way to estimate the *value* of an idea or concept or proposition or doctrine?
- (3) What ought to be understood by the term Truth?

It is no part of my present purpose to consider the worth of Pragmatism's answers to those basic questions. That is a task for you. I intend merely to point out the connection of Pragmatism with Charles S. Peirce.

Peirce has been rightly called a "seminal thinker." A seminal thinker is one whose thinking sometimes produces what we may call *seed*-thoughts, thoughts, that is, that germinate and grow and develop into light-giving theories or doctrines, thus advancing human knowledge, insight, and understanding.

Pragmatism sprang from a seed-thought planted by Peirce. That is well known and acknowledged by all. The seed was planted by him as a maxim of logic. Here it is, twice stated in his own words:

"Consider what effects, that might conceivably have practical bearings, we conceive the object of our con-

CHARLES SANDERS PEIRCE AS A PIONEER

ception to have. Then, our conception of these effects is the whole of our conception of the object."

"In order to ascertain the meaning of an intellectual conception one should consider what practical consequences might conceivably result by necessity from the truth of that conception; and the sum of these consequences will constitute the entire meaning of the conception."

That seed was planted in 1878. It seems pretty dry. As a matter of fact it lay dormant for 20 years. Imagine Peirce's surprise and his astonishment when it then burst forth as a vital, throbbing, world-shaking philosophy, sponsored, elaborated, ardently advocated by such men as William James, John Dewey, F.C.S.Schiller, G.Papini, and others.

Was Peirce quite pleased with their popular expositions of his idea? If you think so, read his *Seven Lectures on Pragmatism* (*Collected Papers*, Vol.V), especially the beginning of the first one. Peirce thought these eloquent gentlemen were a bit too glib, too "lively," too enthusiastic, too superficial. By the "new pragmatists" and by literary folk his idea of pragmatism had been, in his opinion, so distorted that in 1905 he abandoned the term *Pragmatism* and adopted for his own use in its stead the term *Pragmaticism*. Note the delicious humor in his announcement of the change:

"So then, the writer, finding his bantling 'pragmatism' so promoted, feels that it is time to kiss his child good-by and relinquish it to its higher destiny; while to serve the precise purpose of expressing the original definition, he begs to announce the birth of

the word 'pragmaticism,' which is ugly enough to be safe from kidnappers."

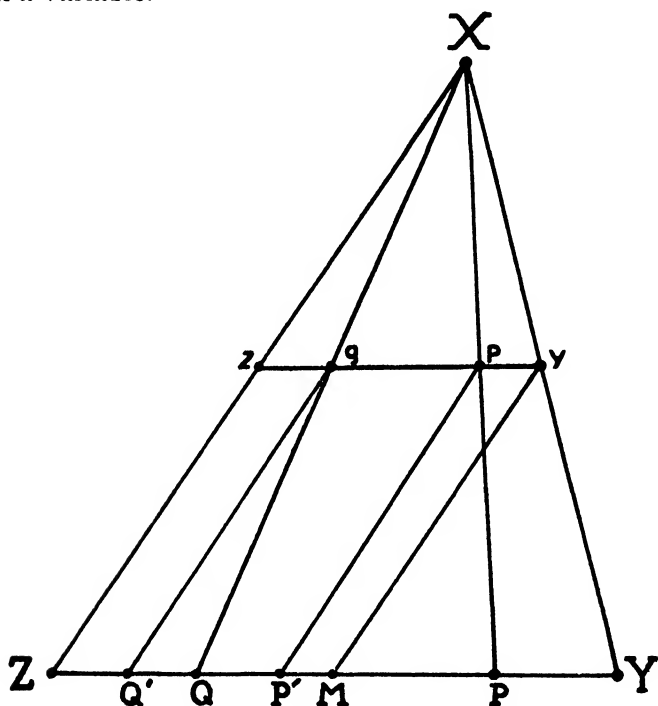
You may be specially interested to consider whether James's conception of Truth can be reconciled with Peirce's conception of it. If you are, read what James says about Truth in his above-cited book and then read in volume V of the *Collected Papers* the following passages: pp.242-3, 252-3, 268, 279, 392-4.

Infinite

You probably know that in mathematics the term infinite is employed in two widely different senses. Let us briefly recall them to mind. One of them has to do with certain increasing variables; the other, with a certain kind of aggregates, or sets, or collections, or ensembles, or classes of things (elements or items of some kind, no matter what kind).

If a variable V be such that in course of its variation it becomes and remains greater than any prescribed amount, however large, we say that V becomes *infinite*, or approaches ∞ (infinity). For a simple example, let V be the fraction $1/D$, where we will suppose D to be an infinitesimal—a variable having zero for its limit. Now choose any finite magnitude M , however great. As D approaches its limit, the variable $1/D$ will plainly become and remain greater than M . To describe this behavior of the variable we say that it becomes infinite or approaches ∞ . In such cases the symbol ∞ does not denote a definite finite quantity which V approaches as a *limit*. In the case supposed, V has no limit. The phrase "becomes infinite" or the equivalent phrase "approaches ∞ " is only a short way of saying that V becomes greater than any prescribed

finite quantity and remains greater through all succeeding stages of V 's variation. That is the meaning, the sole meaning, of the elliptic phrase. The sense explained is the older and more familiar sense of the word infinite. It may be called the *dynamic* sense because it is always connected with a variable.



The newer and less familiar sense of the adjective may be called the *static* sense because it always refers to a certain kind of class, and a class is not a variable, but is a constant, fixed, static affair. The static sense of the word infinite may be clearly defined as follows: If a class C (of items, called its members) be such that its members can be paired in a one-to-one way with a part of its members

(constituting a subclass C'), the class C is said to be an *infinite* class. For a simple example, let C be the class of positive integers, and let C' be the class of even integers. Each integer may be paired with its double. Thus the members of C get paired with those of C' , a part of C . Obviously the pairs are: 1, 2; 2, 4; 3, 6; 4, 8; ; n , $2n$; Hence the integers constitute an infinite class.

For a geometric example, let C be the class of all the points of the segment YZ and let C' be the class of the points of the segment MZ , the half of YZ . Construct the figure as indicated. The lines from X pair the points of YZ with those of yz ; that is, Y with y , P with p , and so on. The parallels from the points of yz pair these with the points of MZ ; that is, y with M , p with P' , and so on. Now, as P corresponds to p in the first pairing, and p to P' in the second pairing, we agree to let P correspond to P' in a third pairing. Thus the points of YZ are paired with those of MZ . So we see that the points of any given segment constitute an infinite class.

But is it legitimate to say that such infinite classes exist? Is it legitimate to speak of "all" of the integers or of "all" the points of a line? Most mathematicians have said yes; some have said no. I am not about to debate the question here, but, without asserting the existence of such classes, I am going to speak of them *as if* they exist, thus using the phrase, as if, in the sense explained by Vaihinger in his *Philosophie des Als Ob*. So doing, we may say that the mathematical theory of infinite aggregates is a branch of mathematics dealing with a certain "As If" variety of subject-matter. The literature of this branch is, as you doubtless know, vast, exceedingly fascinating, and in large part seemingly indispensable. Many of its theorems are, quite literally, stunning. For example, the integers are

exactly as numerous as the integers and fractions taken together. Again, a sphere surface of a microscopically small radius contains precisely as many points as one of radius equal to a million light years. To see it, let the surfaces be concentric, and reflect that any radial line pierces the little surface in a point P and the big one in a point Q. Thus the points of the little surface are paired with those of the huge one. Outstanding among the originators and builders of this towering mathematical edifice are Galileo, Bernhard Bolzano, Richard Dedekind, and especially Georg Cantor.

Peirce has not been generally credited with having been a pioneer in this field but he was. To define an infinite class as being one that is not finite is well enough provided we know what a finite class is. Peirce began by defining the finite. Primarily a logician, he was interested to delimit the applicability of various forms of argument. His definition of the finite is—and the fact is notable—in terms of such a form. More than fifty years ago, in volume IV of the *American Journal of Mathematics*, he defined a *finite* class to be one to which the syllogism of transposed quantity is applicable. An example of such a syllogism is this:

Every Hottentot kills a Hottentot,
No Hottentot is killed by more than one
Hottentot
(Hence)
Every Hottentot is killed by a Hottentot.

If the argument is valid the class of Hottentots is finite, and conversely, (by Peirce's definition); if invalid, the class is infinite (in accord with the foregoing definition).

Peirce alleges that this and other ideas of his were ap-

propriated by Dedekind in his *Was Sind und Was Sollen die Zahlen* without giving him credit (*Collected Papers*, Vol.IV, p.268; Vol.V, p.110). Peirce's discussions of infinity are numerous and some of them extensive. In my opinion an excellent subject for an essay for the degree of master of arts in philosophy or mathematics would be The Contributions of Charles Sanders Peirce to the Theory of Infinite Aggregates.

You are doubtless aware that some notion of Infinity has figured conspicuously in philosophic, theological and scientific thought from the earliest times. A Critical History of the Concept of Infinity is a genuine desideratum, and would, I think, be an admirable dissertation for the degree of doctor either in Philosophy or in Mathematics. In Lecture XV of my *Mathematical Philosophy* I have endeavored to show what one chapter of such a critique might be.

Propositional Function

The term *propositional function* was coined by Bertrand Russell for use in the philosophy of mathematics. The meaning of the term may be best shown by a few simple examples. Let me say first that I am going to use the term proposition to denote a statement that is true or else false, not both true and false, and not neither true nor false.

Consider the statements

- (1) Mahatma Ghandi is a man,
- (2) x is a man,
- (3) $2 + 5 = 9$,
- (4) $2 + y = 9$.

Obviously (1) and (3) are propositions, one true, the other

false. Just as obviously (2) and (4), though each has the form of a proposition, are not propositions, being neither true nor false, since they say nothing definite. They are propositional functions. The x and the y are indeterminate, they are *variables*, being ready to take almost any meanings you please to give them. Such assignable meanings are their values. Substitution of values for the variables yields propositions. These are values of the propositional functions, called *functions* because they are themselves variables depending on the variables in them, and called *propositional* functions because their values are propositions.

A propositional function may contain any number of variables, as seen in such examples as: X marries Y ; $x + 2y - 4z = 0$; A ran against B for the office C of the city D in the campaign E . The variables need not be denoted by x , y , etc. They may be ordinary dictionary words.

Consider the familiar geometric postulate: Two points determine a line. If the terms point and line have not been given definite meanings, they are variables, and the statement is, not a proposition, but a propositional function. It might as well or better be replaced by the statement, $2x$'s determine a y . Thus you see that the postulates of a branch of mathematics are propositional functions. I will leave it to you to show that the theorems are also propositional functions. The one set of functions implies the other set. If the postulates are P and one of the theorems is T , then the statement, P implies T , is neither a postulate nor a theorem; it is not even a propositional function but is a (mathematical) proposition, being quite independent of the variables, or of any meanings that may be assigned to the variables, in P and T .

Thus you begin to see what the rôle of propositional functions is and how very important it is. A little reflection will suffice to show you that most of the so-called propositions in the books of the world are not genuine propositions, since their terms are indeterminate, but are propositional functions and are, accordingly, neither true nor false. No wonder that the centuries of human history have been filled with disputation.

I have said that the term propositional function was invented by Bertrand Russell. But the idea denoted by it was not originated by him. It was familiar to Charles Sanders Peirce and employed by him before Russell was born. But Peirce did not call it a propositional function, he called it a *rheme* or *rhema*, (plural) *rhemata*. Let the following words of Peirce bear witness:

If parts of a proposition be erased so as to leave *blanks* in their places, and if these blanks are of such a nature that if each of them be filled by a proper name the result will be a proposition, then the blank form of proposition which was first produced by the erasures is termed a *rheme*. According as the number of blanks in a rheme is 0, 1, 2, 3, etc., it may be termed a *medad*, *monad*, *dyad*, *triad*, etc., rheme (Collected Papers, Vol.II, p.155).

If in any written statement we put dashes in place of two or more demonstratives or pro-demonstratives, the professedly incomplete representation resulting may be termed a *relative rheme* . . . All rhemata higher than the dual (dyadic) may be considered as belonging to one and the same order; and we may say that all rhemata are either singular, dual, or plural Such, at least, is the doctrine I have

been teaching for twenty-five years, and which, if deeply pondered, will be found to enwrap an entire philosophy (Vol.III, 262-263).

The last of the foregoing statements was made by Peirce 67 years ago.

Paradoxes

In recent years philosophic students of mathematics and logic have been a good deal concerned with certain curious phenomena which they call "paradoxes."

When I was a boy I was carefully taught that a paradox is a proposition that seems obviously false but is nevertheless true. In this sense of the term, the proposition that there are just as many whole numbers as there are whole numbers and fractions together is a paradox. So, too, is the proposition that there are precisely as many points inside a small circle or sphere as inside a large one. And there are hosts of other examples. This is the sense in which the term is employed by Bernhard Bolzano in his book *Paradoxien des Unendlichen*. But the meaning of the term as used by the students above referred to is very different.

In their sense of the term a paradox is a statement which presents itself in the guise of a perfectly respectable proposition but which, when scrutinized, is found to have the diabolic character of being *false if true*, and *true if false*. Perhaps the most familiar example is the paradox of the village barber *who shaves all and only the men in the town who do not shave themselves*. Does he shave himself? It is plain that, if he does, he does not, and that, if he does not, he does. Another specimen is afforded by the assertion that there is a class C whose members are all the classes such that no one of them is a member of itself. Is the class C a member of itself? A little reflection suffices to show

that, if it is, it is not, and that, if it is not, it is. In volume I of their *Principia Mathematica* Whitehead and Russell have listed and discussed many such paradoxes.

It is not difficult to manufacture them. If we deliberately make them as we make playthings and use them as such, they are sufficiently amusing. But when, as often happens, they suddenly and unexpectedly spring up in the course of our soberest thinking to baffle progress and compel retreat, they are far from amusing. As I have said, we can make paradoxes, but can we *avoid* making them? If not, we have to wonder whether even our proudest discourse, which we boast of as logically rigorous, is not in fact rendered quite silly by the lurking presence in it of some stultifying and mocking paradox.

Long before *Principia Mathematica* appeared numerous paradoxes were noted and acutely discussed by Peirce. But he did not call them paradoxes.* He called them *Insolubilia* (Collected Papers, Vol.II, pp.200,370; Vol.III, p.281; Vol.IV, pp.52-54; Vol.V, pp.203-212).

The following shows how Peirce handled one of them:

“Given the following proposition:

This assertion is not true:

is that assertion, which proclaims its own falsity, and nothing else, true or false? Suppose it true. Then,

Whatever is asserted in it is true,

But that it is not true is asserted in it;

∴ By Barbara, That it is not true is true;

∴ It is not true.

Besides, if it is true, that it is true is true. Hence,

That it is not true is not true,

CHARLES SANDERS PEIRCE AS A PIONEER

But that it is not true is asserted in the proposition;

∴ By Darapti, Something asserted in the proposition is not true:

∴ The proposition is not true.

On the other hand, suppose it is not true. In that case,

That it is not true is true,

But all that the proposition asserts is that it is not true;

∴ By Barbara, All that the proposition asserts is true;

∴ The proposition is true.

Besides, in this case,

Something the proposition asserts is not true,

But all that the proposition asserts is that it is not true;

∴ By Bokardo, That it is not true is not altogether true;

∴ That it is true is true;

∴ It is true."

Relations

Another field in which Peirce did extensive pioneering work is the immense field of Relations.

Most words, perhaps all of them, denote relations either directly or indirectly. Familiar specimens are: father, mother, child, brother, sister, uncle, king, subject, citizen, inhabitant, above, below, in, out, greater, less, better, worse, major, taller, shorter, area, diameter, derivative, integral, volume, lower, sweeter, savior, president, premier, and so on and on.

Relations are so numerous and so various that most of them have no names. For it seems that everything, great or small, is related to everything else. Can you think of something that has no relations at all? Can you think of two things absolutely unrelated to each other? Our lives are immersed in a sea of relations whether we work or play or rest or sleep or dream or wake. When we study Science, said Henri Poincaré, we are studying, not things, but their relations. The universe appears to be an infinite net-work of relations. Being, said Lotze, consists of relations. To be is to be related. Not to be related is not to be.

Relations, we say, relate *things*. Are things, too, relations? If so, there is nothing but relations. That view is advocated by Charles W. Wood in his recent book entitled *The Myth of the Individual*. His contention, stated in my own words, is, in effect, that the web of reality is woven of relations and that the so-called things or individuals which the relations are said to connect are themselves composed of relations, being, so to say, only branch-points or nodes or knots or ganglia formed in the web by the meeting and intersection of relations.

Wood's idea, which seems a bit mystical, may not be quite defensible. Nevertheless relations are so essential in our world, so omnipresent that one can not help wondering why it never occurred to anyone, until quite recently, to make Relations as such a subject of scientific research, so that relation-theory would become a great branch of Science. It is, of course, true, always has been, and always will be, that whenever we are investigating an object we are studying that object's relations to other things. But such incidental, only half-conscious, and ad hoc study of special relations is far from being the same thing as the deliberate and conscious investigation of relations as such

and in general, with a view to establishing a Logic of Relations.

Aristotle's Logic is a great work, but it is mainly concerned with *classes* and *propositions*, only slightly with *relations*. The three things—classes, propositions, and relations—though they are radically different, are so interconnected that one cannot think much about or discuss any one of them without, in some measure, thinking about and discussing the other two. But it did not occur to Aristotle that, in slighting the subject of relations, he was slighting the most momentous part of the proper subject-matter of logic. And nothing was done to correct the great error for over 2000 years.

It is to the genius of Augustus DeMorgan (1806–1871) that we owe the first weighty contribution to the theory of relations. I refer to his memoir "On the Syllogism IV, and the Logic of Relations" (*Cambridge Philosophical Transactions*, vol.10).

Peirce became acquainted with that work of DeMorgan's in 1866. Prior to that date he had independently discovered the potentiality and desirability of such a logic and in subsequent years he brought it to the estate of a well established science, and this was prior to Russell's *Principles of Mathematics* and Whitehead and Russell's *Principia Mathematica*, in which works relation-theory plays a major rôle.

In the study of relations one finds on the very threshold of the subject a large number of pretty obvious, very interesting, and very important facts. Let us glance at a few of them.

We agree to say that a thing x has a given relation R to a thing y by writing xRy . Here x is a *referent*, and y a *relatum*. The referents and relata of a relation are its

terms. Note that a relation has a direction, so to speak, or *sense*, running from referent to relatum. If xRy , then y has to x a relation called the *converse* of R and denoted by \bar{R} .

A relation has two aspects: an *intensional* aspect, as when we say, for example, that "the relation father means a man who has begotten a child"; and *extensional* aspect, as when we say that the relation of father *is* the *class* of pairs or *couples*, (x_1, y_1) , (x_2, y_2) , and so on, where the x 's are fathers and the y 's are children.

Regarding a relation as a class of couples, we can say that two relations do or do not *intersect* according as the two classes have or do not have one or more couples in common. Also we can say that one relation is part of another, if all the couples of the former belong to the latter. And so on. It is found that, for scientific purposes, the *extensional* aspect of relations is more important than the *intensional* aspect.

The referents of R constitute its *domain*; the relata, its *co-domain*; the referents and relata together, its *field*. The domain and co-domain may have no term in common, as in case of the relation, husband; or they may overlap, without coinciding, as in case of the relation, brother; or they may coincide, as in case of the relation, spouse.

The domain and co-domain of a relation may both be *finite* classes, or both *infinite* classes, or either may be finite and the other infinite. I leave it to you to exemplify these cases.

Relations are endlessly diversified but, happily, they fall into classes, some very large, some smaller. Let us note a few of the big classes.

A relation may be *dyadic* or *triadic*, etc. It is dyadic if defined or definable by a propositional function having two and but two variables, as x is president of y ; triadic if

the defining function has three variables, as x sold y to z ; and so on. I have spoken and shall continue to speak of dyadics only.

A relation may be *symmetric* or *asymmetric* or *non-symmetric*. R is symmetric if xRy implies yRx , as in case of the relation, equals, or spouse. R is asymmetric if xRy implies that yRx is false, or impossible; an example is the relation, greater than, or mother. R is called non-symmetric if, when xRy , we sometimes have yRx , and sometimes not; an example is the relation, brother, or lover.

Again, a relation may be *transitive* or *intransitive* or *non-transitive*. R is transitive if xRy and yRz together imply xRz , as in case of the relation, equals, or greater, or less, or like. R is intransitive if xRy and yRz together imply that xRz is false, as in case of the relation, parent, or uncle. R is non-transitive, if, when xRy and yRz , it sometimes happens that xRz and sometimes not; an example is the relation, unequals, or hater.

I strongly recommend as a fine exercise the finding of relations that exemplify the various foregoing distinctions and combinations of them. Find, for example, a relation that is both symmetric and transitive.

Relations can be operated on in various ways. One of the most important operations is called *relative multiplication*. If R and S be two relations, the co-domain of R and the domain of S may or may not have one or more terms in common. I will leave it to you to find an R and an S such that, if y be a relatum of R it can not be a referent of S . If R and S be such that there is a y for which we have xRy and ySz , then the consequent relation of x to z is called the *relative product* of R and S (in the order named) and is denoted by the symbol $R | S$, and we have $xR | Sz$. If, for example, R be father, and S be mother,

then $R | S$ is maternal grandfather. Obviously relative multiplication is not commutative.

At present the Logic of Relations, or as Peirce called it, the Logic of Relatives, contains an immense number of theorems. I will close this little talk about relations by giving one of the theorems, one whose sweep is sublime. It is this:

$$\text{Conv.} \sim (R | S) = \check{S} | \check{R};$$

that is, the converse of the relative product of R and S is the relative product of the converse of S by the converse of R .

It is easy to prove as follows. Let R' stand for Conv. of $(R | S)$. We have to prove that, if $xR'z$ then $x\check{S} | \check{R}z$, and, conversely, that, if $x\check{S} | \check{R}z$, then $xR'z$. Since $xR'z$, then $zR | Sx$; hence there is a y such that zRy and ySx ; hence $y\check{R}z$ and $x\check{S}y$; hence $x\check{S} | \check{R}z$. Conversely, if $x\check{S} | \check{R}z$, then there is a y for which $x\check{S}y$ and $y\check{R}z$; hence ySx and zRy ; hence $zR | Sx$, and hence $xR'z$.

A final remark: In my article entitled "A Glance at Some of the Ideas of Charles Sanders Peirce" (*Scripta Mathematica*, vol.III, No.1, January, 1935) I have signalized a number of other respects in which Peirce did important work as a creator or pioneer and have indicated, by citing volume and page, where the work occurs in the *Collected Papers*.

THE FOURTH DIMENSION
AND RELATIVITY

By LEOPOLD INFELD

THE FOURTH DIMENSION AND RELATIVITY¹

THE problem of the fourth dimension seems to the layman both obscure and mysterious. He is not even sure that the mathematicians and physicists who claim to understand this problem really do so. The statement that the world is a four dimensional curved space earns either blind admiration for the imaginative physicists or, more often, results in the conviction that modern physics must be absurd since it contradicts common sense.

I shall now try to tell you something about the problem of the fourth dimension and its connection with relativity. You can hardly expect me to give a complete explanation of this subject. Perhaps I may succeed, in the short time at my disposal, in convincing you that there is not the slightest trace of mystery in the problem. Should you feel that you need the thrill connected with a mystery to help you admire the achievements of the scientists, then my talk will sorely disappoint you. Every great scientific idea can be clearly, unambiguously formulated. Science has no use for mystery and there should be no place for it in scientific popularization.

But let us return to our problem of the fourth dimension. I should like to present a dialogue with a clever pupil asking me just the right questions! Let us begin.

P. It must be wonderful for the mathematicians and physicists to be able to imagine the fourth dimension. Could I learn how it is done?

¹ A Radio Broadcast given over WNYG, May 7, 1938, under the auspices of the Galois Institute of Mathematics, Long Island University.

GALOIS LECTURES

I. A physicist can imagine the fourth dimension as little or as much as you can.

P. That is hard for me to believe.

I. I shall try to convince you. Let me ask you some questions and you will answer them as well as you can and then ask me further questions. Let us begin with a simple one, forgetting for the moment the problem of the fourth dimension. Can you imagine a one dimensional space?

P. I certainly can. Just as well as I can imagine a line. The line and the one dimensional space are, for me, one and the same thing.

I. Let us pause for a moment to introduce a new word. You said that, for you, a line represents a one dimensional space. In ordinary language space always means a three dimensional space, so in order to avoid a muddle, let us use the word *continuum*. If you accept my proposition then I may say that for you a line is a one dimensional continuum. I repeated your statement, changing only one word, that is "space" into "continuum." I hope you do not object! But what do you mean by the sentence that a line is a one dimensional continuum? Suppose that I have no intuitive notion of space or even if I have I do not wish to use it. How would you then explain to me your sentence about the line representing a one dimensional continuum?

P. Let me think for a moment. Perhaps this would do. A point which remains on a line can move only in one of two opposite directions, in the same way as a car on a very narrow highway. Is not this restriction of motion characteristic of the one dimensional space?

I. Your explanation does not sound very encouraging although I think I see your point and understand the intuitive idea you are trying to express. But I could get you into trouble by asking what you mean by a point

THE FOURTH DIMENSION AND RELATIVITY

moving along a line in the same or in opposite directions. You are using in your explanation very complicated ideas indeed: those of motion and direction! The mathematician understands the sentence about the line representing a one dimensional continuum in a much simpler way which I shall try to explain. The explanation may bore you a little but once you grasp the idea the rest will be easy. Take two arbitrary points on the line. Call one of them *A* and the other *B*. Or, simpler still, denote one of them by the number "zero" and the other by the number "one." Between the points "0" and "1" there are an unlimited number of points. If we want to denote, without repetition, every point between "0" and "1" by a number in this interval then, we must use all the numbers between "0" and "1." In this way two different points will always be denoted by two different numbers and vice versa. Every number between "0" and "1" will denote a point between *A* and *B*. More generally speaking we can say what we understand by the sentence, "all points on a line form a one dimensional continuum." This sentence means: a definite number corresponds to every point; a definite point corresponds to every number.

P. But if this is so then a time or temperature scale also forms one dimensional continuum because a definite number corresponds to each instant of time or to each temperature. And on the other hand a definite instant and a definite temperature correspond to each number.

I. This is true. We can equally well represent a one dimensional continuum by all temperatures or by all instants or by all points on a line. In each of these cases and in many others the characteristics of a continuum are preserved.

P. I believe you. What you really did was to substitute

the word "continuum" for the word "space" and then settle its meaning in such a way that not only a line but also many physical quantities which can be pictured as lines and are therefore associated with all numbers, may be regarded as a continuum. I take this for granted but should like to know where it leads. What do you understand by a two dimensional continuum? A plane?

I. A plane is a good example. I can give you many others. Imagine that our earth is a perfect sphere. Each point on the earth can be characterized by longitude and latitude. Every schoolboy denotes in this way the situation of different towns. More abstractly we can say: a pair of numbers denoting the longitude and latitude corresponds to every point on the earth. Every pair of numbers within certain limits can be characterized by a point. The continuum is not one dimensional but two dimensional because not one number but a pair of numbers corresponds to a point. Or take another example. Imagine that you, or preferably someone else, is lying sick in a hospital. A most conscientious nurse takes the patient's temperature, say, every hour. The result can be drawn on a chart in the form of a curve such as you must certainly have seen. The measurements result in a curve in a two dimensional plane representing the two dimensional temperature-time continuum. Two lines with one common point determine a plane just as two one dimensional continua of temperature and time determine a two dimensional time-temperature continuum. We can be unnecessarily learned and say that the nurse taking her patient's temperature registers points in this two dimensional continuum. Every point is, in fact, characterized by two numbers: one denoting time, the other temperature.

P. This means then that in certain cases we can form a

THE FOURTH DIMENSION AND RELATIVITY

two dimensional continuum from two one dimensional continua, such as a plane from two intersecting lines. Or, in other words, we can compose a higher dimensional continuum from the one dimensional ones.

I. If you have once understood this, neither 4 nor 4000 dimensions can frighten you. Instead of 4000 numbers you can say a point in a 4000 dimensional space. Instead of all possible sets of 4000 numbers you can say all points in a 4000 dimensional space, or, still simpler, the whole 4000 dimensional space. This is merely a convenient way of expressing what you mean. When speaking of the simplest case, the one dimensional continuum, we can picture it as a line, the two dimensional continuum as a surface and the three dimensional continuum as ordinary space. We cannot visualize a 4 or a 4000 dimensional continuum but this need not bother us. The most important thing is that we know what we are talking about and we do already know that points in a many dimensional space must be regarded as sets of as many numbers as the space has dimensions. It is not nonsense to say "points in a four dimensional space," for we are well aware what this means. The question whether we are able to visualize something concretely is of no importance in pure mathematics. Once we have understood that a point in space means a set of numbers, we are unlikely to be scared by a great number of dimensions.

P. I agree with you that, in this way, a mathematician may please himself with as many dimensions as he likes. But you cannot deny that we live in a three dimensional space. Or, using your terminology, I could say that in our space the positions of points are characterized by three numbers. To characterize the position of a particular corner of the table I must know three things, its distance from

the ceiling and from two perpendicular walls. A stone thrown out of the window moves in this three dimensional space. Let us assume for the moment that you have really convinced me that I cannot forbid the mathematician to speak of any number of dimensions he likes. I believe, nevertheless, that I can still limit the number of dimensions about which the physicist may speak. We live in a three dimensional space and that is that. We observe all the events in nature in this three dimensional space. Does not the introduction of a 4-dimensional space into physics mean bringing in metaphysics, which is outside our experience?

I. If we approach this point slowly I hope I shall be able to show you that, in talking of a four dimensional world, the physicist is not playing at metaphysics. We already know that three numbers correspond to any point in space and that a point in space corresponds to any three numbers. It is characteristic of the continuum that we can find a point as near as we wish to any given point. But all this is scarcely physics. In physics we must investigate the world of events. To state just when and where Brutus stabbed Caesar we must use four numbers. Three of these characterize the point in space and the fourth characterizes the time. In the expression, "it happened in *Rome*," Rome stands for three space numbers. If we read that a train arrives at Washington at 8:30 we actually note the characterization of an event where, "railway station in Washington," means a place ordinarily and much more accurately characterized by three numbers. To characterize the sending or the reception of a light signal we must use four numbers. Three of them state the place and the fourth the instant at which the event occurs. If this is so, then you understand that the world of events forms a four dimen-

THE FOURTH DIMENSION AND RELATIVITY

sional continuum. Four numbers correspond to every event and an event corresponds to every four numbers.

P. Already I must raise some objections. Are you allowed to combine the three dimensional space and the one dimensional time continuum into an artificial four dimensional time-space continuum? Are not these two continua of quite different natures?

I. Try to be consistent. Why did you not object to the nurse plotting the graph of her patient's temperature by combining time and temperature into a two dimensional time-temperature continuum? What about the many graphs in your daily paper where, for example, the rise or fall of trade forms one continuum and time forms the other?

P. I see your point of view but my objections still hold. I will try to formulate them. For example, in driving a car I pass, at one moment, a tree, at the next a street corner and a moment later a traffic circle. I should say that the car changed its position with time. The picture which I use is *dynamic*. I observe changes of position in our three dimensional space, changes going on in time. My description has nothing to do with four dimensions. For me time is something quite different from space. I call my picture *dynamic* because the positions of the body, that is the car, change with time. You would describe the motion in a different way. You would say: If the car passes a tree the event is characterized by four numbers, three denoting the position, the fourth the instant. According to you, therefore, we have a point in a four dimensional time-space continuum. For me it is motion in space, a change of position with time. For you the picture is different, it is *static*. You picture the events as stationary points placed in a four dimensional space-time continuum. I agree that your picture of motion is quite consistent but it seems to me very

artificial and, in any case, I see no advantage to be gained by using it.

I. But you agree that I may use my static picture if I wish. That, after all, is something. I must admit that motion can be pictured either your way or mine. Either dynamically against the background of a three dimensional space or statically against the background of a four dimensional space. The two pictures are equivalent and you, therefore, see no reason for choosing the static. You believe, if I may say so, that it shows bad taste to use the more artificial, static four dimensional picture, rather than the simple dynamic picture of events changing with time in a three dimensional space. But at least you must agree with me that there is nothing mysterious in all this.

P. Yes, I must agree and for the simple reason that I understand your point. I do not think that I could understand if there were any mystery about it. But one thing still escapes me. What has all this to do with relativity? You have not yet used the word relativity, but the concept of a four dimensional world is, in most peoples' minds, connected with the theory of relativity.

I. That is a different story. I am afraid that my explanation will not satisfy you. After all, to understand the theory of relativity demands a special study and you can scarcely expect me to explain the structure of the theory in a few words. The only thing I can do is to formulate dogmatically some statements concerning the four dimensional time-space structure and its rôle in the theory of relativity.

Remember that from the point of view of classical physics I may use either of two pictures of motion, the dynamic or the static, the three or the four dimensional background. The theory of relativity, however, taught us

THE FOURTH DIMENSION AND RELATIVITY

that these two pictures are not equivalent but that the static or four dimensional is preferable. The aim of every physical theory is to formulate physical laws simply and logically. To understand and order our experiences we must use the four dimensional background. According to the theory of relativity the four dimensional time-space continuum and not the three dimensional continuum should be regarded as the scene of events in our physical world. The static and not the dynamic picture of motion is the objective one. Once again, as so often in physics, we are forced to abandon our intuitive view and retreat to a position which although further from our immediate experience enables us to understand more fully the world of our impressions.

There is one point more. We have already quoted examples of two dimensional continua, the plane and the sphere. These two continua have different geometrical properties, though both represent two dimensional continua. The statement that the events in our world are best described against the background of a four dimensional time-space continuum expresses little. What we still need is some information as to the kind of continuum and its geometrical properties. The answer to this question would lead us directly to the problems of relativity. Let us, however, stop here and summarize briefly what we have learned.

We call a set of n numbers a point in an n dimensional continuum. A line is a one dimensional, a surface a two dimensional, and our physical space a three dimensional continuum. We can speak of continua of more dimensions without caring whether or not we can imagine them. Every physical event can be described as a point in a four dimensional continuum. This is as true for classical physics

as for the theory of relativity. But the theory of relativity has shown the static picture of motion as something existing in a four dimensional continuum as more objective than the dynamic picture of motion as something happening in a three dimensional continuum. The four dimensional time-space continuum is the scene of events of our physical world and it is the aim of the theory of relativity to describe its geometrical properties.

